# FGAHOI: Fine-Grained Anchors for Human-Object Interaction Detection

Shuailei Ma, Yuefeng Wang, Shanze Wang, and Ying Wei

**Abstract**—Human-Object Interaction (HOI), as an important problem in computer vision, requires locating the human-object pair and identifying the interactive relationships between them. The HOI instance has a greater span in spatial, scale, and task than the individual object instance, making its detection more susceptible to noisy backgrounds. To alleviate the disturbance of noisy backgrounds on HOI detection, it is necessary to consider the input image information to generate fine-grained anchors which are then leveraged to guide the detection of HOI instances. However, it is challenging for the following reasons. $i$) how to extract pivotal features from the images with complex background information is still an open question. $ii$) how to semantically align the extracted features and query embeddings is also a difficult issue. In this paper, a novel end-to-end transformer-based framework (FGAHOI) is proposed to alleviate the above problems. FGAHOI comprises three dedicated components namely, **multi-scale sampling (MSS)**, **hierarchical spatial-aware merging (HSAM)** and **task-aware merging mechanism (TAM)**. MSS extracts features of humans, objects and interaction areas from noisy backgrounds for HOI instances of various scales. HSAM and TAM semantically align and merge the extracted features and query embeddings in the hierarchical spatial and task perspectives in turn. In the meanwhile, a novel training strategy **Stage-wise Training Strategy** is designed to reduce the training pressure caused by overly complex tasks done by FGAHOI. In addition, we propose two ways to measure the difficulty of HOI detection and a novel dataset, $i.e.$, HOI-SDC for the two challenges (**Uneven Distributed Area in Human-Object Pairs** and **Long Distance Visual Modeling of Human-Object Pairs**) of HOI instances detection. Experiments are conducted on three benchmarks: HICO-DET, HOI-SDC and V-COCO. Our model outperforms the state-of-the-art HOI detection methods, and the extensive ablations reveal the merits of our proposed contribution. The code is available at https://github.com/xiaomabufei/FGAHOI.

**Index Terms**—Human-Object Interaction, FGAHOI, Fine-Grained Anchors, Noisy Background, Semantically Aligning.

✦

## 1 INTRODUCTION

HUMAN-Object interaction (HOI) detection, as a downstream task of object detection [1], [2], [3], [4], [5], has recently received increasing attention due to its great application potential. For successful HOI detection, it needs to have the ability to understand human activities which are abstracted as a set of <human, object, action> triplets in this task, requiring a much deeper understanding for the semantic information of visual scenes. Without HOI detection, machines can only interpret images as collections of object bounding boxes, i.e., AI systems can only pick up information such as 'A man is on the bike' or 'A bike is in the corner', but not 'A man rides a bike'.

Spanning the past and the present, the existing HOI detection approaches [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21] tend to fall into two categories, namely two-stage and one-stage methods. Conventional two-stage methods [7], [8], [10], [12], [13], [14], [18], [20], [22], [23], [24], [25], as an intuitive approach, detect human and object instances by leveraging the off-the-
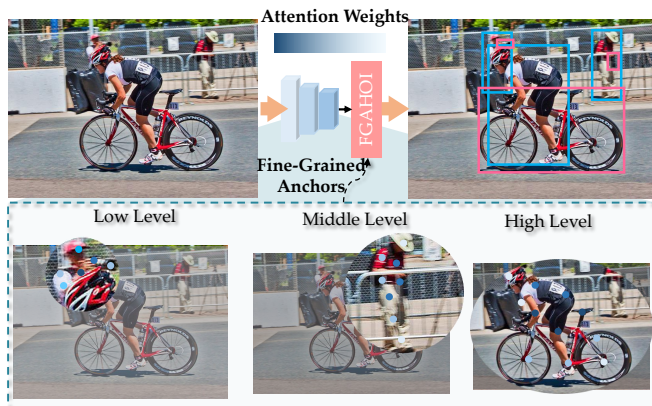


Fig. 1: FGAHOI leverages the query embeddings and multi-scale features to generate fine-grained anchors and the corresponding weights for HOI instances of diverse scales. Then, they guide the decoder to aid key semantic information of HOI instances to the content embeddings and translate the content embeddings to HOI embeddings for predicting all elements of the HOI instances.

shelf object detector [1], [3], [4], utilizing the visual features extracted from the located areas to recognize action classes. To fully leverage the visual features, several methods [7], [10], [14], [20], [22], [23], [24], [25] separately extract visual features of human-object pairs and spatial information from the located area in a multi-stream architecture, fusing them in a post-fusion strategy. In the meanwhile, several approaches [8], [10], [20], [23], [24] employ the existing pose

- Shuailei Ma, Yuefeng Wang are with College of Information Science and Engineering, Northeastern University, Shenyang, China, 110819. E-mail: {xiaomabufei, wangyuefeng0203} @gmail.com
- Shanze Wang is with Changsha Hisense Intelligent System Research Institute Co., Ltd. and Information Technology R&D Innovation Center of Peking University, Shaoxing, China. E-mail: szgg0099@gmail.com
- Ying Wei is the corresponding author, with College of Information Science and Engineering, Northeastern University, Shenyang, China, 110819. E-mail: weiying@ise.neu.edu.cn

estimation methods, such as [26], [27], [28] to extract pose information and fuse it with other features to predict the action class. In addition, some works [8], [12], [13], [18], [29] leverage the graph neural network to extract complex semantic relationship between humans and objects. However, the difficulties encountered in the two-stage approach lie mainly in the effective fusion of human-object pairs and complex semantic information. Besides, owing to the limitations of the fixed detector and some other components (pose estimation etc.), the two-stage method can only achieve a sub-optimal solution.

To achieve high efficiency, one-stage approaches [6], [9], [11], [15], [17], [21], [30], [31] which utilize interaction points between the human-object pairs to simultaneously predict human and object offset vectors and action classes, are proposed to detect human-object pairs and recognize interactive relationships in parallel. However, when the human and object in the image are far apart from each other, these methods are disturbed by ambiguous semantic features. The one-stage methods do not achieve much attention until the appearance of the Detection Transformer (DETR) [32] and QPIC [19] applies it for HOI detection. Then, plenty of transformer-based works [6], [9], [16], [17], [33] attempt to solve the HOI detection with different encoder-decoder structures and backbone models.

In comparison to object instances, HOI instances have a greater span of spatial, scale and task. In most HOI instances, there is a certain distance between human and objects and their scale varies enormously. Compared with simple object classification, it is necessary to consider more information between human-object pairs rather than the features of humans and objects for interaction classification. Therefore, the detection is more susceptible to distractions from noisy backgrounds. However, most recent works [19], [33] use object detection frameworks [32], [34] directly for HOI detection by simply adding the interaction classification head, ignoring these problems. Inspired by [34] which leverages the reference points to guide the decoding process, we propose to leverage fine-grained anchors to guide the detection of HOI instances and protect it from noisy backgrounds. To generate fine-grained anchors for kinds of HOI instances, it is obviously necessary to consider the input image features. There are, however, two inevitable challenges that arise as a result of this. *i*) it is difficult to extract pivotal features from the images which contain noisy background information. *ii*) how to semantically align and merge the extracted features with query embeddings is also an open question.

In this paper, we propose a novel transformer-based model for HOI detection, i.e., FGAHOI: Fine-Grained Anchors for Human-Object Interaction Detection (as shown in Fig.1). FGAHOI leverages the **multi-scale sampling mechanism (MSS)** to extract pivotal features from images with noisy background information for variable HOI instances. Based on the sampling strategy and initial anchor generated by the corresponding query embedding, MSS could extract hierarchical spatial features of human, object and the interaction region for each HOI instance. Besides, the **hierarchical spatial-aware (HSAM)** and **task-aware merging mechanism (TAM)** are utilized to semantically align and merge the extracted features with the query embeddings.

HSAM merges the extracted features in the hierarchical spatial perspective according to the cross-attention between the features and the query embeddings. Meanwhile, the extracted features are aligned towards the query embeddings, according to the cross-attention weights of the merging process. Thereafter, TAM leverages the switches which dynamically switch ON and OFF to merge the input features and query embeddings in the task perspective.

According to experiment results, we investigate that it is difficult of the end-to-end training approach to allow the transformer-based models to achieve optimal performance when more complex task requirements are required. Inspired by the stage-wise training [35], [36] for LTR [37], we propose a novel stage-wise training strategy for FGAHOI. During the training process, we add the important components of the model in turn to clarify the training direction of the model at each stage, so as to maximize the savings in the training cost of the model.

To the best of our knowledge, there are no measurements for the difficulty of detecting HOI instances. We investigate that two difficulties lie in the detection of human-object pairs, *i.e.*, **Uneven Distributed Area in Human-Object Pairs** and **Long Distance Visual Modeling of Human-Object Pairs**. In this paper, we propose two measurements and a novel dataset (HOI-SDC) for these two challenges. HOI-SDC eliminates the influence of other factors (Too few training samples of some HOI categories, too tricky interaction actions, et.al.) on the model training and focuses on the model for these two difficult challenges. Our contributions can be summarized fourfold:

- We propose a novel transformer-based human-object interaction detector (FGAHOI) which leverages input features to generate fine-grained anchors for protecting the detection of HOI instances from noisy backgrounds.
- We propose a novel training strategy where each component of the model is trained in turn to clarify the training direction at each stage, in order to maximize the training cost savings.
- We propose two ways to measure the difficulty of HOI detection and a dataset, *i.e.*, HOI-SDC for the two challenges (Uneven Distributed Area in Human-Object Pairs and Long Distance Visual Modeling of Human-Object Pairs) of detecting HOI instances.
- Our extensive experiments on three benchmarks: HICO-DET [38], HOI-SDC and V-COCO [39], demonstrate the effectiveness of the proposed FGAHOI. Specifically, FGAHOI outperforms all existing state-of-the-art methods by a large margin.

## 2 RELATED WORKS

**Two-stage HOI Detection Approaches**: The two-stage HOI detection approaches [7], [8], [10], [12], [13], [14], [18], [20], [22], [23], [24], [25], [29] employ the off-the-shelf object detector [1], [3], [4] to localize humans and objects. Afterwards, the features of backbone networks inside the human and objects regions are cropped. Part of the two-stage methods [8], [12], [13], [18], [29] treat the human and objects feature as nodes and employ graph neural networks [40]
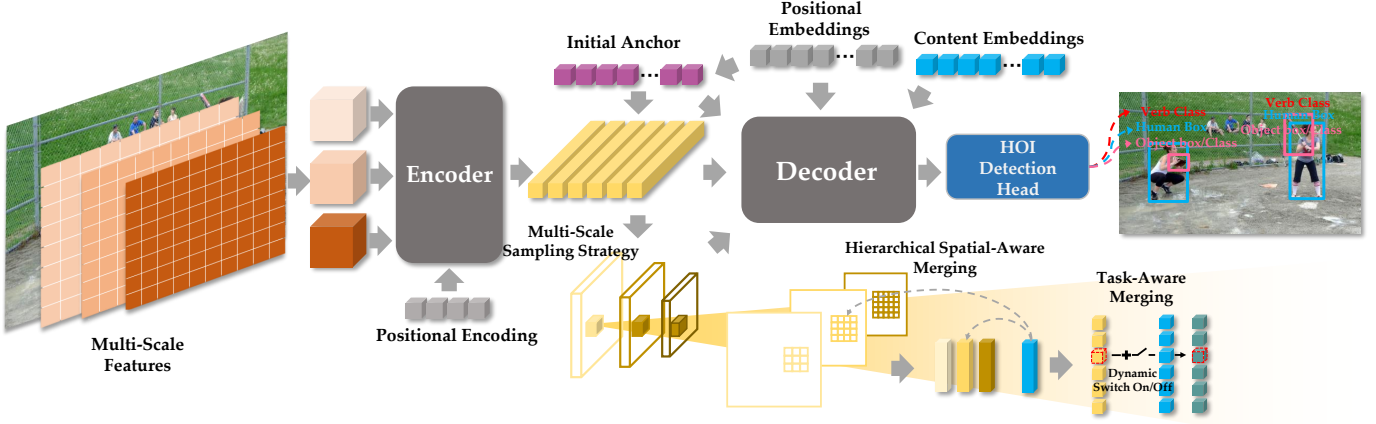
Fig. 2: This figure illustrates the overall structure of FGAHOI. FGAHOI utilizes a hierarchical backbone and a deformable encoder to extract the semantic features in a multi-scale approach. In the decoding phrase, FGAHOI leverages the **multi-scale sampling**, **hierarchical spatial-aware merging** and **task-aware merging** mechanism to align input features with query embeddings and assist the generation of fine-grained anchors for the translation of HOI embeddings. At the back end of the pipeline, HOI detector leverages the HOI embeddings and initial anchor to predict all elements of the HOI instances.

to predict action classes. The other part of the two-stage approach [7], [10], [14], [20], [22], [23], [24], [25] leverages multi-stream networks to extract diverse information from cropped regions, such as human features, object features, spatial information and human pose information. Then, the information is fused to predict the action in a post-fusion strategy. Two-stage methods mainly concentrate on predicting the action class in the second stage. Nevertheless, the quality of cropped features from the first stage cannot be guaranteed in most cases, so the method cannot achieve an optimal solution. More importantly, integrating semantic information of human-object pairs requires massive time and computing resources.

**One-stage HOI Detection Approaches:** The traditional one-stage approaches [9], [11], [15], [31] use interaction points or union regions to detect human-object pairs and identify interactive action classes in parallel. However, these methods which and are hampered by distant human-object pairs, require a gathering and pairing process. With the creation of DETR [32], one-stage approaches have become the current mainstream. QPIC [19] converts the object detection head of DETR into an interaction detection head to predict HOI instance directly. HOITrans [17] combines transformer [41] and CNN [42] to straightly predict HOI instances from the query embeddings. AS-Net [6] and HOTR [9] each propose a two-branch transformer method that consists of an instance decoder and an interaction decoder to predict the boxes and action classes in parallel. CDN [16] proposes a cascade disentangling decoder to decode action classes. QAHOI [33] directly combines Swin Transformer [43] and deformable DETR [34] to predict HOI instances.

**Anchor-Based Object Detection Transformer:** Deformable DETR [34] first introduces the reference point concept, where the sampling offset is predicted by each reference point to perform deformable cross-attention. To facilitate extreme region discrimination, Conditional DETR [44] reformulates the attention operation and rebuilt positional queries based on reference points. Anchor DETR [45] pro-

poses to explicitly capitalize on the spatial prior during cross-attention and box regression by utilizing a predefined 2D anchor point $[cx, cy]$. DAB-DETR [46] extends such a 2D concept to a 4D anchor box $[cx, cy, w, h]$ and proposed to refine it layer-by-layer. SAM-DETR [47] proposes directly updating content embeddings by extracting salient points from image features. In this paper, we propose a novel decoding process for HOI detection. The alignment and fine-grained anchor generation is proposed to align the multi-scale features with HOI query embeddings and generate fine-grained anchors for the diverse HOI instances with variable spatial distribution, scales and tasks. Then, the fine-grained anchors guide the deformable attention process in aiding key information to query embeddings from noisy backgrounds.

## 3 PROPOSED METHOD

In Sec.3.1, we show the overall architecture of FGAHOI. Then, we describe the multi-scale feature extractor in Sec.3.2. We introduce the multi-scale sampling strategy in Sec.3.3.1. The hierarchical spatial-aware, task-aware merging mechanism and the decoding process is proposed in Sec.3.3.2, Sec.3.3.3 and Sec.3.3.4, respectively. In Sec.3.4, we present the architecture of the HOI detection head. In Sec.3.5, the stage-wise training strategy, loss calculation and inference process is illustrated.

### 3.1 Overall Architecture

The overall architecture of our proposed FGAHOI is illustrated in Fig 2. For a given image $x \in \mathbb{R}^{H \times W \times 3}$, FGAHOI firstly uses a hierarchical backbone network to extract the multi-scale features $Z_i \in \mathbb{R}^{\frac{H}{4 \times 2^i} \times \frac{W}{4 \times 2^i} \times 2^i C_s}, i = 1, 2, 3$. The multi-scale features are then projected from dimension $C_s$ to dimension $C_d$ by using 1×1 convolution. After being flattened out, the multi-scale features are concatenated to $N_s$ vectors with $C_d$ dimensions. Afterwards, along with

supplementary positional encoding $p \in \mathbb{R}^{N_s \times C_d}$, the multi-scale features are sent into the deformable transformer encoder which consists of a set of stacked deformable encoder layers to encode semantic features. The encoded semantic features $M \in \mathbb{R}^{N_s \times C_d}$ are then acquired. In the decoding process, the content $C$ and positional $P$ embeddings are both a set of learnable vectors $\{v_i \mid v_i \in \mathbb{R}^{c_d}\}_{i=1}^{N_q}$. The positional embeddings $P$ first generate the initial anchor $A \in \mathbb{R}^{N_q \times 2}$ according to a linear layer. The positional $P$, content $C$ embeddings, inital anchor $A$ and encoded features $M$ are simultaneously sent into the decoder $F_{decoder}(\cdot, \cdot, \cdot, \cdot)$ which is a set of stacked decoder layers. In every decoder layer, the initial anchor first leverages the multi-scale sampling strategy to sample the multi-scale features corresponding to the content embeddings. The sampled features assist the generation of fine-grained anchors and corresponding attention weights through the hierarchical spatial-aware and task-aware merging mechanism. The HOI embeddings $H = \{h_i \mid h_i \in \mathbb{R}^{c_d}\}_{i=1}^{N_q}$ are translated from the query embeddings $Q$ through the fine-grained anchors, attention weights and the deformable attention. The HOI embeddings $H$ are acquired as $H = F_{decoder}(M, P, C, A)$. Eventually, the HOI detector leverages the HOI embeddings $H$ and initial anchor to predict the HOI instances $< b_h, b_o, c_o, c_v >$, where $b_h$, $b_o$, $c_o$ and $c_v$ stands for the human box coordinate $(x, y, w, h)$, object box coordinate, object class and verb class, respectively.

## 3.2 Multi-Scale Features Extractor

High-quality visual features are a prerequisite for successful HOI detection. For extracting the multi-scale features with long-range semantic information, FGAHOI leverages the multi-scale feature extractor which consists of a hierarchical backbone network and a deformable transformer encoder to extract features, the folumation is as Equation.1:

$$M = F_{encoder}(F_{flatten}(\phi(x)), p, s, r, l) \in \mathbb{R}^{N_s \times C_d}, \quad (1)$$

where $F_{encoder}(\cdot)$, $F_{flatten}(\cdot)$ and $\phi(\cdot)$ denotes the encoder, flatten operation and backbone network, respectively. $p$ is the position encoding, $s$ is the spatial shape of the multi-scale features, $r$ stands for the valid ratios and $l$ represents the level index corresponding the multi-scale features. The hierarchical backbone network is flexible and can be composed of any convolutional neural network [42], [48], [49], [50] and transformer backbone network [43], [51], [52], [53], [54], [55], [56], [57]. However, CNN is poor at capturing non-local semantic features like the relationships between humans and objects. In this paper, we mainly use Swin Transformer tiny and large version [43] to enhance the ability of feature extractor for extracting long-range features.

## 3.3 Why FGAHOI Decodes Better?

During the decoding process, the fine-grained anchors can be regarded as a positional prior to let decoder focus on the region of interest, directly guiding the decoder to aid semantic information to the content embeddings which are used to predict all elements of the HOI instances. Therefore, fine-grained anchors play the following two crucial roles in HOI detection. *i)* Fine-grained anchors directly determine

whether the information gained from input features to content embeddings is instance-critical or noisy background information. *ii)* Fine-grained anchors determine the quality of alignment between the query embeddings and multi-scale features of input scenarios. Both are crucial factors for the quality of decoding results. The existing methods [33], [34] directly utilize the query embedding to generate fine-grained anchors based on the initial anchor, without considering the multi-scale features of the input scenarios and the semantic alignment between the query embedding and the input features at all. Our FGAHOI proposes a novel fine-grained anchors generator which consists of **multi-scale sampling**, **hierarchical spatial-aware merging** and **task-aware merging mechanism** (as shown in Fig.3). The generator adequately leverages the initial anchor, multi-scale features and query embeddings for generating suitable fine-grained anchors for diverse input scenarios and aligning semantic information between different input scenarios and query embeddings. The formulation of FGAHOI decoding process is as follows:

$$H = \text{Defattn}(\text{Task}(\text{Hier Spatial}(\{x_s^i\}, C_u), C_u), M, C_u), \quad (2)$$

where $C_u$ is the content embeddings updated by the positional embeddings, Defattn represents the deformable attention, $x_s^i$ represents the sampled features of the $i$-th level features. $M$ is the encoded input features.

### 3.3.1 Multi-Scale Sampling Mechanism

The HOI instances contained in the input scenarios usually vary in size, where some instances taking up most of the area in the input scenarios and others occupying perhaps only a few pixels. Our FGAHOI aims at detecting all instances in the scene, regardless of the size. Therefore, when using the initial anchor to sample the multi-scale features, for shallow features mainly used to detect instances of small size, the sampling strategy only samples a small range of features around the initial anchor. In contrast, for deep features mainly used to detect instances of large size, the sampling strategy samples a large range of features around the initial anchor. As shown in Fig.3 (b), in the generator, the encoded features are first reshaped to the original shape. Based on the initial anchor, generator leverages the sampling strategy to sample multi-scale features as follows:

$$x_s^i = F_{sample}( reshape(M)^i, A, size^i, bilinear ), \quad (3)$$

where $size^i$ ($i = 0, 1, 2$) denotes the sampling size of the $i$-th level features. $M$ is the encoded input features. $A$ is the initial anchor. Inspired by [58], we utilize bilinear interpolation in the sampling strategy.

### 3.3.2 Hierarchical Spatial-Aware Merging Mechanism

In order to better utilize the hierarchical spatial information of sampled features for aligning content embeddings with the sampled features, we propose a novel hierarchical spatial-aware merging mechanism (HSAM) which utilizes the content embeddings to extract hierarchical spatial information and merge the sampled features, as shown in Fig.3 (c). The content embeddings are first updated by the positional embeddings and multi-head self-attention mechanism as follows:

$$C_u = C + F_{\text{MHA}}\left((C + P)W^q, (C + P)W^k, CW^v\right), \quad (4)$$
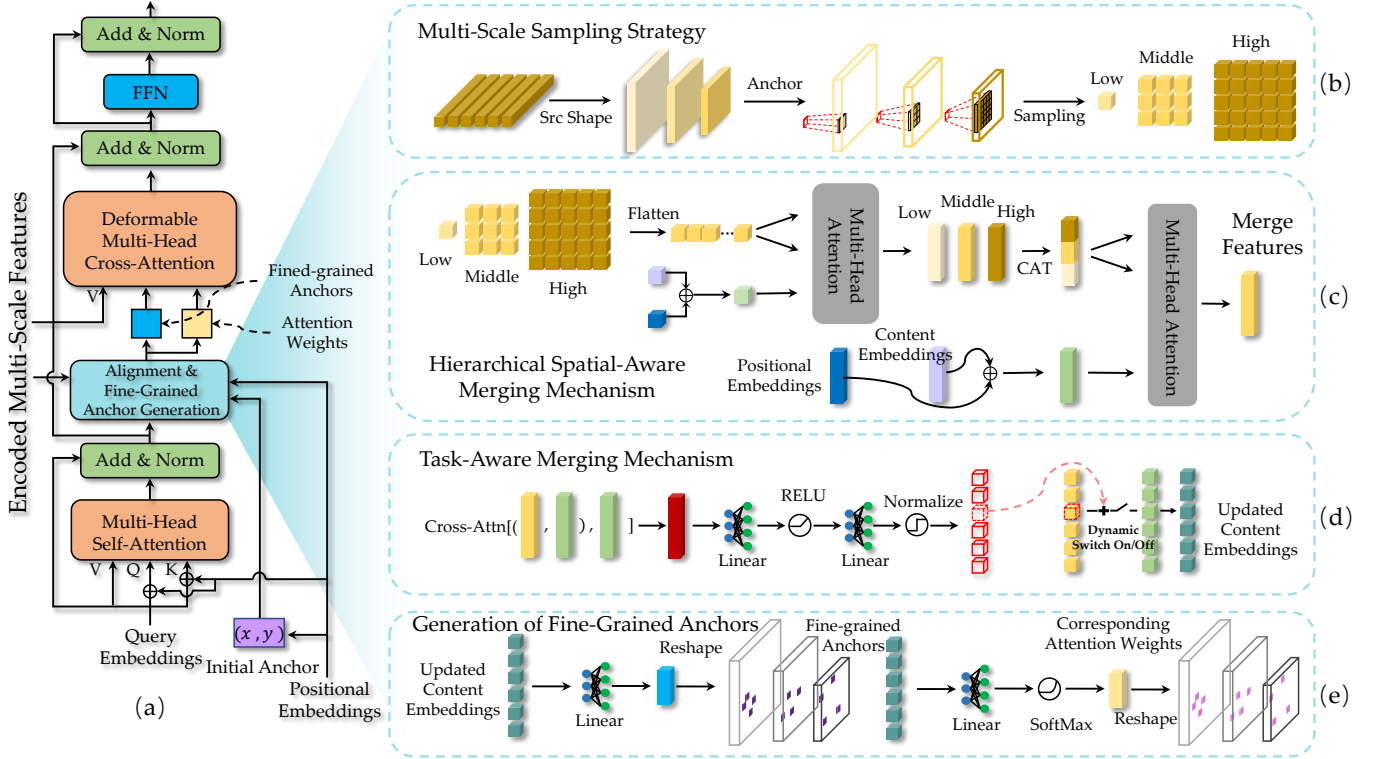
Fig. 3: The architecture of FGAHOI's decoder. (a) Illustration of FGAHOI's decoding process. (b) Illustration of Multi-scale sampling mechanism. (c) Illustration of Hierarchical spatial-aware merging mechanism. (d) Illustration of Task-aware merging mechanism. (e) Generation process of fine-grained anchors and the corresponding attention weights.

where $W^q$, $W^k$ and $W^v$ denotes the parameter matrices for query, key and value in the self-attention mechanism, respectively. $F_{\text{MHA}}(\cdot)$ is the multi-head attention mechanism. $C$ and $P$ represents the content and position embeddings, respectively. Then, the updated content embeddings are leveraged to merge the sampled features, the formulation is as follows:

$$x_m^i = F_{\text{concat}}\left(\text{head}_1, \ldots, \text{head}_{\text{N}_\text{H}}\right) W^O,$$
$$\text{where} \quad \text{head}_\text{n} = \text{Softmax}\left(\frac{(C_u W_\text{n}^q) \cdot (x_s^i W_\text{n}^k)^T}{\sqrt{d_k}}\right)(x_s^i W_\text{n}^v). \quad (5)$$

Where $x_m^i$ represents the merged features of the $i$-th level sampled features. $C_u$ is the content embeddings updated by the positional embeddings. $W^O$ denotes the parameter matrices for multi-head concatenation. $W_n^q$, $W_n^k$ and $W_n^v$ denote the parameter matrices for query, key and value of n-th attention head. $F_{\text{concat}}$ is the concatenating operation. $d_k = \frac{N_{hd}}{N_H}$, $N_{hd}$ is the hidden dimensions, and $N_H$ is the number of attention head.

Following the merging of the sampled features at each scale based on spatial information, the merged features at each scale are first concatenated together as follows:

$$X_m = F_{\text{concat}}(\{x_m^i\}_{i=0,1,2}) \in \mathbb{R}^{B \times N_q \times N_L \times N_{hd}}, \quad (6)$$

where $N_L$ is the number of multi-scale, $x_m^i$ represents the merged features of the $i$-th level sampled features, $X_m$ is the concatenated multi-scale features and merged by the scale-aware merging mechanism as follows:

$$X_u = F_{\text{concat}}\left(\text{head}_1, \ldots, \text{head}_\text{h}\right) W^O,$$
$$\text{where} \quad \text{head}_\text{n} = \text{Softmax}\left(\frac{(C_u W_\text{n}^q) \cdot (X_m W_\text{n}^k)^T}{\sqrt{d_k}}\right)(X_m W_\text{n}^v). \quad (7)$$

Where $X_u$ is the merged multi-scale features for updating the content embeddings.

### 3.3.3 Task-Aware Merging Mechanism

Considering diverse HOI instances, the task-aware merging mechanism is proposed to fuse the merged multi-scale features and content embeddings and align the content embeddings with the merged feature in the task-aware perspective, as shown in Fig.3 (e). It leverages the merged multi-scale features and content embeddings to generate dynamic switch for selecting suitable channel in the merging process. Content embedding and multi-scale information after fusion are first stitched together, the formulation is as follows:

$$X = F_{stack}(C_u, X_u) \in \mathbb{R}^{B \times N_q \times (2 \times N_{hd})}. \quad (8)$$

Where $C_u$ is the content embeddings updated by the positional embeddings, $X_u$ is the merged multi-scale features. Thereafter, we use cross-attention mechanism to update these as follows:

$$X_{switch} = F_{\text{concat}}\left(\text{head}_1, \ldots, \text{head}_\text{h}\right) W^O,$$
$$\text{where} \quad \text{head}_\text{n} = \text{Softmax}\left(\frac{(C_u W_\text{n}^q) \cdot (X W_\text{n}^k)^T}{\sqrt{d_k}}\right)(X W_\text{n}^v). \quad (9)$$

TABLE 1: Instance statistics of two difficulties. We quantify all the instances in the HAKE-HOI [20] dataset according to two newly proposed metrics and divide them into ten intervals.

| Dataset | IMI | $\text{IMI}_0$ | $\text{IMI}_1$ | $\text{IMI}_2$ | $\text{IMI}_3$ | $\text{IMI}_4$ | $\text{IMI}_5$ | $\text{IMI}_6$ | $\text{IMI}_7$ | $\text{IMI}_8$ | $\text{IMI}_9$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| HAKE-HOI | $\text{num}_{AR}$ | 104243 | 65499 | 44303 | 31241 | 21982 | 11888 | 4670 | 1818 | 598 | 168 |
| | $\text{num}_{LR}$ | 424 | 1243 | 1784 | 3043 | 8668 | 70191 | 83314 | 79427 | 34017 | 4299 |
| SDC_Train | $\text{num}_{AR}$ | 62526 | 30235 | 16346 | 12013 | 10269 | 11189 | 4223 | 1540 | 423 | 139 |
| | $\text{num}_{LR}$ | 177 | 515 | 874 | 1656 | 5208 | 48798 | 38517 | 29544 | 20265 | 3349 |
| SDC_Test | $\text{num}_{AR}$ | 24737 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | $\text{num}_{LR}$ | 153 | 415 | 464 | 834 | 2704 | 20167 | 0 | 0 | 0 | 0 |

Then, the generated information is utilized to gain the dynamic switch for merging, the formulation is as follows:

$$Switch^{\gamma} = F_{normalize}(F_{mlp}(X_{switch}))^{\gamma} \in \mathbb{R}^{B \times N_q \times 2 \times 2}, \quad (10)$$

where $Switch^{\gamma}$ is the dynamic switch for $\gamma$-th dimension of the merged features. $F_{hsigmoid}(\cdot)$ and $F_{mlp}(\cdot)$ denote the hard sigmoid and feed forward network which consists of two linear layers and one Relu activation layer, respectively. Inspired by [59], the merging mechanism is designed as follows:

$$U^{\gamma} = F_{Max}\{Switch^{\gamma}_{i,0} \odot X_u^{\gamma} + Switch^{\gamma}_{i,1}\}_{i=0,1} + C_u^{\gamma}, \quad (11)$$

where $U^{\gamma}$ is $\gamma$-th features of content embeddings updated by the merged multi-scale features. $F_{Max}$ is the max operation.
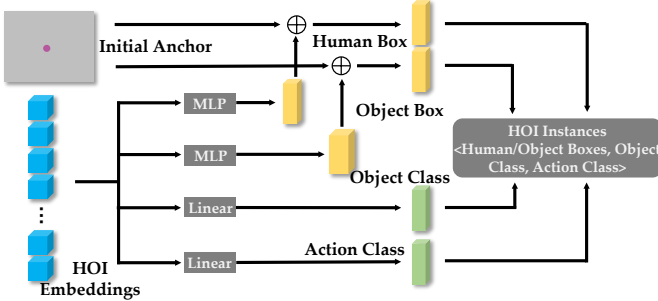


Fig. 4: The prediction process of the HOI detection head. See sec 3.4 for more details.

### 3.3.4 Decoding with Fine-Grained Anchor

As shown in Fig.3 (e), the updated content embeddings are used to generate fine-grained anchors and attention weights. According to the linear layer, reshape operation and softmax function, the formulation is as follows:

$$\mathcal{A} = F_{lin\&res}(U) \in \mathbb{R}^{B \times N_q \times N_H \times N_L \times N_{\mathcal{A}} \times 2}, \quad (12)$$

$$\mathcal{W} = F_{lin\&res\&soft}(U) \in \mathbb{R}^{B \times N_q \times N_H \times N_L \times N_{\mathcal{A}}}, \quad (13)$$

As shown in Fig.3 (a), the fine-grained anchors and attention weights are utilized to aid semantic features from the encoded features of the input scenarios to the content embeddings, the formulation is as follows:

$$\mathcal{P}_q = \sum_{n=1}^{N_H} W_n \left[ \sum_{l=1}^{N_L} \sum_{k=1}^{N_{\mathcal{A}}} \mathcal{W}^l_{nqk} \cdot W'_n x^l \left( \mathcal{A}^l_{nqk} \right) \right], \quad (14)$$

where $\mathcal{P}_q$ is the extracted semantic information used for translating $q$-th content to HOI embeddings. $\mathcal{A}^l_{nqk}$ and $\mathcal{W}^l_{nqk}$ represent the $k$-th fine-grained anchors and corresponding attention weights of the $n$-th attention head for the $q$-th query embedding. Both $W_n$ and $W'_n$ are parameter matrices of the $n$-th attention head. $N_{\mathcal{A}}$ is the number of fine-grained anchors of each scale in one attention head.

### 3.4 HOI Detection Head

FGAHOI leverages a simple HOI detection head to predict all elements of HOI instances. As shown in Fig.4, the detection head utilizes the HOI embeddings and the initial anchor to localize the human and object boxes. In this process, each initial anchor acts as the base point for the bounding boxes of the corresponding pair of a human and an object, the formulation is as follows:

$$b_h = F_{mlp}(H)[\cdots,:2] + initial\ anchor \quad \in \mathbb{R}^{N_q \times 4}, \quad (15)$$

$$b_o = F_{mlp}(H)[\cdots,:2] + initial\ anchor \quad \in \mathbb{R}^{N_q \times 4}, \quad (16)$$

$$c_o = F_{linear}(H) \quad \in \mathbb{R}^{N_q \times num_o}, \quad (17)$$

$$c_v = F_{linear}(H) \quad \in \mathbb{R}^{N_q \times num_v}, \quad (18)$$

where $F_{mlp}$ denotes the feed forward network consists of three linear layers and three relu activation layers. $F_{linear}$ stands for the linear layer. $num_o$ and $num_v$ are the number of object and action classes, respectively. $H$ denotes the HOI embeddings.

### 3.5 Training and Inference

#### 3.5.1 Stage-wise Training

Inspired by the stage-wise training approach [35], [36] which decouples feature learning and classifier learning into two independent stages for LTR [37], we propose a novel stage-wise training strategy for FGAHOI. We start by training the base network (FGAHOI without any merging mechanism) in an end-to-end manner. We then add the merging mechanism in turn to the trained base network for another short period of training. In this phrase, the parameters of the trained base network are leveraged as pretrained parameters and no parameters are fixed during the training process.

Fig. 5: Visualization of HOI detection. Humans and objects are represented by pink and blue bounding boxes respectively, and interactions are marked by grey lines linking the box centers. Kindly refer to Sec. 5.6.1 for more details.
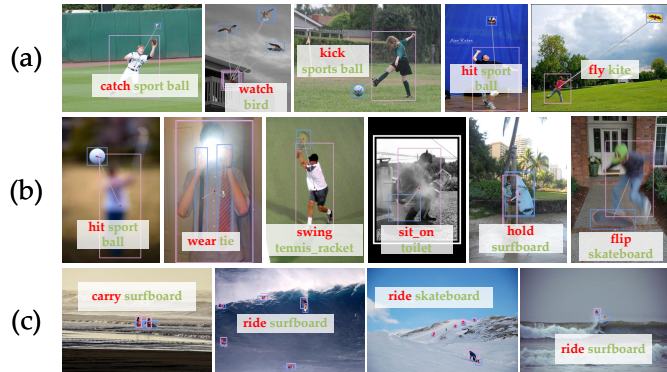


Fig. 6: (a) illustrates the excellent long-range visual modelling capabilities. (b) demonstrates remarkable robustness. (c) shows the superior capabilities for identifying small HOI instances. Kindly refer to Sec. 5.6.1 for more details.

### 3.5.2 Loss Calculation

Inspired by the set-based training process of HOI-Trans [17], QPIC [19], CDN [16] and QAHOI [33], we first use the bipartite matching with the Hungarian algorithm to match each ground truth with its best-matching prediction. For subsequent back-propagation, a loss is then established between the matched predictions and the matching ground truths. The folumation is as follows:

$$L = \lambda_o L_c^o + \lambda_v L_c^v + \sum_{k \in (h,o)} \left( \lambda_b L_b^k + \lambda_{GIoU} L_{GIoU}^k \right), \quad (19)$$

where $L_c^o$ and $L_c^v$ represent the object class and action class loss, respectively. We utilize the modified focal loss function [60] and sigmoid focal loss function [61] for $L_c^v$ and $L_c^o$, respectively. $L_b$ is the box regression loss and consists of the $L1$ Loss. $L_{GIoU}$ denotes the intersection-over-union loss, the same as the function in QPIC [19]. $\lambda_o$, $\lambda_v$, $\lambda_b$ and $\lambda_{GIoU}$ are the hyper parameters for adjusting the weights of each loss.

### 3.5.3 Inference

The inference process is to composite the output of the HOI detection head to form HOI triplets. Formally, the $i$-th output prediction is generated as $< b_i^h, b_i^o, argmax_k c_i^{hoi}(k) >$. The HOI triplet score $c_i^{hoi}$ combined by the scores of action $c_i^v$ and object $c_i^o$ classification, formularized as $c_i^{hoi} = c_i^v \cdot c_i^o$.

## 4 PROPOSED DATASET

There are two main difficulties existing with human-object pairs. $i$) Uneven size distribution of human and objects in human-object pairs. $ii$) Excessive distance between person and object in human-object pairs. To the best of our knowledge, there are no relevant metrics to measure these two difficulties. In this paper, we propose two metrics $AR$ and $LR$ for measuring these two difficulties. Then two novel challenges corresponding to these two difficulties are proposed. In addition, we propose a novel **S**et for these **D**ouble **C**hallenges (HOI-SDC). The data is selected from HAKE-HOI [20] which is re-split from HAKE [62] and provides 110K+ images. HAKE-HOI has 117 action classes, 80 object classes and 520 HOI categories.

Fig. 7: Comparison of fine-grained anchors between FGAHOI and QAHOI. We visualize the fine-grained anchors corresponding to all attention heads and the corresponding attention weights, where the shades of colors correspond to the magnitude of the weights. Obviously, FGAHOI is more accurate in focusing on humans, objects and interaction areas. Kindly refer to Sec. 5.6.2 for more details.

## 4.1 HOI-UDA

We propose a novel measurement for the challenge of **Uneven Distributed Area in Human-Object Pairs**, the formulation is as follow:

$$AR = \frac{Area_h \cdot Area_o}{Area_{hoi}^2}, \tag{20}$$

where $Area_h$, $Area_o$ and $Area_{hoi}$ denote the area of human, object and HOI instances, respectively (as shown in Fig.8 (a)). We quantify all the instances in the HAKE-HOI into ten intervals and count the number of instances of each interval in the second and fifth row of Table.1. To better evaluate the ability of the model to detect HOI for human-object pairs with uneven distributed areas, we specially select 24737 HOI instances of $IMI_0^{UDA}$ in testing set.

## 4.2 HOI-LDVM

A novel measurement for the challenge of **Long Distance Visual Modeling of Human-Object Pairs** is proposed in Eq.21.

$$LR = \frac{L_h + L_o}{L_{hoi}}, \tag{21}$$

where $L_h$, $L_o$ and $L_{hoi}$ denote the size we define of human, object and HOI instances, respectively (as shown in Fig.8 (b)). The instances are quantified in the third and sixth row of Table.1. To better evaluate the ability of the model to detect HOI for human-object pairs with with long distance, we specially select 24737 HOI instances of $IMI_0^{LDVM} \sim IMI_6^{LDVM}$ in testing set.

## 4.3 HOI-SDC

In order to avoid the training process of the model being influenced by a portion of HOI classes with a very small number of instances, we remove some of the HOI classes containing a very small number of instances and HOI classes with no interaction from the training **S**et for the **D**ouble **C**hallenge. Finally, there are total 321 HOI classes,

TABLE 2: Performance comparison with the state-of-the-art methods on the HICO-DET dataset. 'V', 'S', 'P' and 'L' represent the visual feature, spatial feature, human pose feature and language feature respectively. Fine-tuned Detection means the parameter of the model is pre-trained on the MS-COCO dataset. Backbone with '*' and '+' means that they are pre-trained on ImageNet-22K with 384×384 input resolution. QAHOI(R) represents that the results are reproduced on the same machine with our model. Kindly refer to Sec. 5.4.1 for more details.

| Architecture | Method | Backbone | Fine-tuned | Feature | Default (↑) | | | Known Object (↑) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Full | Rare | Non-Rare | Full | Rare | Non-Rare |
| **Two-Stage Methods** | | | | | | | | | | |
| Multi-stream | No-Frill [23] | ResNet-152 | ✗ | A+S+P | 17.18 | 12.17 | 18.08 | - | - | - |
| | PMFNet [24] | ResNet-50-FPN | ✗ | A+S | 17.46 | 15.65 | 18.00 | 20.34 | 17.47 | 21.20 |
| | ACP [25] | ResNet-101 | ✔ | A+S+L | 21.96 | 16.43 | 23.62 | - | - | - |
| | PD-Net [10] | ResNet-152 | ✗ | A+S+P+L | 22.37 | 17.61 | 23.79 | 26.86 | 21.70 | 28.44 |
| | VCL [7] | ResNet-50 | ✔ | A+S | 23.63 | 17.21 | 25.55 | 25.98 | 19.12 | 28.03 |
| Graph-Based | RPNN [8] | ResNet-50 | ✗ | A+P | 17.35 | 12.78 | 18.71 | - | - | - |
| | VSGNet [13] | ResNet-152 | ✗ | A+S | 19.80 | 16.05 | 20.91 | - | - | - |
| | DRG [12] | ResNet-50-FPN | ✔ | A+S+L | 24.53 | 19.47 | 26.04 | 27.98 | 23.14 | 29.43 |
| | SCG [18] | ResNet-50-FPN | ✔ | A+S | 31.33 | 24.72 | 33.31 | 34.37 | 27.18 | 36.50 |
| **One-Stage Methods** | | | | | | | | | | |
| Interaction points | IP-Net [15] | ResNet-50-FPN | ✗ | A | 19.56 | 12.79 | 21.58 | 22.05 | 15.77 | 23.92 |
| | PPDM [31] | Hourglass-104 | ✔ | A | 21.73 | 13.78 | 24.10 | 24.58 | 16.65 | 26.84 |
| | GGNet [11] | Hourglass-104 | ✔ | A | 23.47 | 16.48 | 25.60 | 27.36 | 20.23 | 29.48 |
| Transformer-Based | HOITrans [17] | ResNet-101 | ✔ | A | 26.60 | 19.15 | 28.54 | 29.1 | 20.98 | 31.57 |
| | HOTR [9] | ResNet-50 | ✗ | A | 23.46 | 16.21 | 25.65 | - | - | - |
| | | ResNet-50 | ✔ | A | 25.10 | 17.34 | 27.42 | - | - | - |
| | AS-Net [6] | ResNet-50 | ✗ | A | 24.40 | 22.39 | 25.01 | 27.41 | 25.44 | 28.00 |
| | | ResNet-50 | ✔ | A | 28.87 | 24.25 | 30.25 | 31.74 | 27.07 | 33.14 |
| | QPIC [19] | ResNet-50 | ✔ | A | 29.07 | 21.85 | 31.23 | 31.68 | 24.14 | 33.93 |
| | | ResNet-50 | ✗ | A | 24.21 | 17.51 | 26.21 | - | - | - |
| | QAHOI [33] | Swin-Tiny | ✗ | A | 28.47 | 22.44 | 30.27 | 30.99 | 24.83 | 32.84 |
| | | Swin-Large$^*_+$ | ✗ | A | 35.78 | 29.80 | 37.56 | 37.59 | 31.66 | 39.36 |
| | QAHOI (R) | Swin-Tiny | ✗ | A | 27.67 | 20.22 | 29.69 | 30.06 | 22.95 | 32.18 |
| | | Swin-Large$^*_+$ | ✗ | A | 35.43 | 29.22 | 37.29 | 37.23 | 31.01 | 39.09 |
| | FGAHOI | Swin-Tiny | ✗ | A | **29.94** | **22.24** | **32.24** | **32.48** | **24.16** | **34.97** |
| | | Swin-Large$^*_+$ | ✗ | A | **37.18** | **30.71** | **39.11** | **38.93** | **31.93** | **41.02** |

74 object classes and 93 action classes. The training and testing set contain 37,155 and 9,666 images, respectively. The detailed distribution of HOI instances is shown in Table.1.



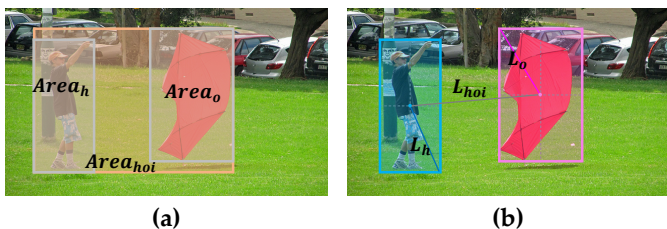**(a)**                                **(b)**

Fig. 8: Proposed metrics for the difficulties existing with HOI instances. (a) Metric for uneven size distribution of humans and objects. (b) Metric for excessive distance between person and object. Kindly refer to Sec. 4.1 and 4.2 for more details.

# 5 EXPERIMENTS

## 5.1 Dataset

Experiments are conducted on three HOI datasets: HICO-DET [38], V-COCO [39] and HOI-SDC dataset

**HICO-DET** [38] has 80 object classes, 117 action classes and 600 HOI classes. HICO-DET offers 47,776 images with

TABLE 3: Performance comparison with the state-of-the-art methods on the HOI-SDC dataset. Kindly refer to Sec. 5.4.2 for more details.

| Dataset | Backbone | Method | mAP$_{role}$ (↑) |
|---|---|---|---|
| HOI-SDC | Swin-Tiny | QAHOI | 19.55 |
| | Swin-Tiny | Baseline | 21.18 |
| | Swin-Tiny | +HSAM | 21.91 |
| | Swin-Tiny | +TAM | 21.84 |
| | Swin-Tiny | FGAHOI | 22.25 |

151,276 HOI instances, including 38,118 images with 117,871 annotated instances of human-object pairs in the training set and 9658 images with 33,405 annotated instances of human-object pairs in the testing set. According to the number of these HOI classes, the 600 HOI classes in the dataset are grouped into three categories: Full (all HOI classes), Rare (138 classes with fewer than ten instances) and Non-Rare (462 classes with more than ten instances). Following HICO [63], we consider two different evaluation settings (the results are shown in Table.2: (1) Known object settings: For each HOI category (such as 'flying a kite'), the detection is only evaluated on the images that contain the target object category (such as 'kite'). The difficulty lies in the local-

TABLE 4: Performance comparison with the state-of-the-art methods on the V-COCO dataset. Kindly refer to Sec. 5.4.3 for more details.

|  | Method | $AP^{S1}_{role}$ ($\uparrow$) | $AP^{S2}_{role}$ ($\uparrow$) |
|---|---|---|---|
| Two-stage Method | VSG-Net | 51.8 | 57.0 |
|  | PD-Net | 52.0 | - |
|  | ACP | 53.2 | - |
| One-stage Method | HOITrans | 52.9 | - |
|  | AS-Net | 53.9 | - |
|  | HOTR | 55.2 | 64.4 |
|  | DIRV | 56.1 | - |
|  | QAHOI(R-50) | 58.2 | 58.7 |
|  | FGAHOI(R-50) | 59.0 | 59.3 |
|  | FGAHOI(Swin-T) | 60.5 | 61.2 |



Fig. 9: The human-object pairs with human overlap $IOU_h$ and object overlap $IOU_o$ both exceeding 0.5 are declared as true positives. Kindly refer to Sec. 5.2 for more details.

ization of HOI (e.g. human-kite pairs) and distinguishing the interaction (e.g. 'flying'). (2) Default setting: For each HOI category, the detection is evaluated on the whole test set, including images containing and without target object categories. This is a more challenging setting because we also need to distinguish background images (such as images without 'kite').

**V-COCO** [39] contains 80 different object classes and 29 action categories and is developed from the MS-COCO dataset, which includes 4,946 images for the test subset, 2,533 images for the train subset and 2,867 images for the validation subset. The objects are divided into two types: "object" and "instrument".

## 5.2 Metric

Following the standard evaluation [21], [39], we use role mean average precious to evaluate the predicted HOI instances. A detected bounding box is considered a true positive for object detection if it overlaps with a ground truth bounding box of the same class with an intersection greater than union ($IOU$) greater than 0.5. In HOI detection, we need to predict human-object pairs. The human-object pairs whose human overlap $IOU_h$ and object overlap $IOU_o$ both exceed 0.5, i.e., min ($IOU_h, IOU_o$) > 0.5 are declared a true positive (as shown in Fig 9). Specifically, for HICO-DET, besides the full set of 600 HOI classes, the role mAP over a rare set of 138 HOI classes that have less than 10 training instances and a non-rare set of the other 462 HOI classes are also reported. Furthermore, we report the role mAP of two scenarios for V-COCO: scenario 1 includes the cases even without any objects (for the four action categories of body motions), while scenario 2 ignores these cases. For HOI-SDC, we report the role mean average precision for the full set of 321 HOI classes.

## 5.3 Implementation Details

The Visual Feature Extractor consists of Swin Transformer and a deformable transformer encoder. For Swin-Tiny and Swin-Large, the dimensions of the feature maps in the first stage are set to $C_s = 96$ and $C_s = 192$, respectively. We pre-train Swin-Tiny on the ImageNet-1k dataset. Swin-Large is first pre-trained on the ImageNet-22k dataset and finetuned
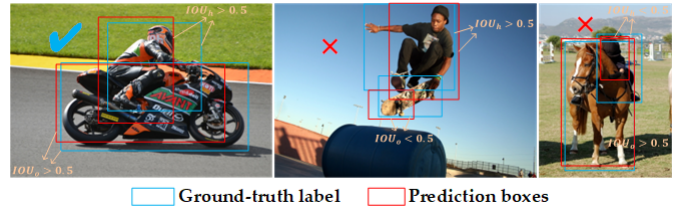
on the ImageNet-1k dataset. Then the weights are used to fine-tune the FGAHOI for the HOI detection task. The number of both encoder and decoder layers are set to 6 ($N_{Layer} = 6$). The number of query embeddings is set to 300 ($N_q = 300$), and the hidden dimension of embeddings in the transformer is set to 256 ($C_d = 256$). In the post-processing phase, the first 100 HOI instances are selected according to object confidence, and we use $\delta$=0.5 to filter the HOI instances by the combined $IOU$. Following Deformable-DETR [34], the AdamW [64] optimizer is used. The learning rates of the extractor and the other components are set to $10^{-5}$ and $10^{-4}$, respectively. We use 8 RTX 3090 to train the model (QAHOI & FGAHOI) with Swin-Tiny. For the model with Swin-Large$^*_+$, we use 16 RTX 3090 to train them. For HICO-DET and HOI-SDC, we train the base network for 150 epochs and carry out the learning rate drop from the 120th epoch at the first stage of training. For subsequent training, we trained the model for 40 epochs, with a learning rate drop at the 15th epoch. For V-COCO dataset, we train the base network for 90 epochs and drop the learning rate from 60th epoch at the first stage of training. For subsequent training, we trained the model for 30 epochs, with a learning rate drop at the 10th epoch.

## 5.4 Comparison with State-of-the-Arts

### 5.4.1 HICO-DET

We compare FGAHOI with the state-of-the-art two-stage and one-stage methods on the HICO-DET dataset and report the results in Table.1. FGAHOI outperforms both state-of-the-art methods. In contrast to the state-of-the-art two-stage method SCG [18], FGAHOI with Swin-Large*+ backbone exceeds an especially significant gain of 5.85 mAP in default full setting, 5.99 mAP in default rare setting, 5.8 mAP in default non-rare setting, 4.56 mAP in known object full setting, 4.75 mAP in known rare settings and 4.52 mAP in known object non-rare setting. For a fair comparison, we used the same machine for the reproduction of the QAHOI (as shown in Table.2 QAHOI(R)). In comparison to the state-of-the-art one-stage method QAHOI, FGAHOI exceeds it in all settings for all backbone networks. For Swin-Tiny backbone network, FGAHOI exceeds an especially significant gain of 2.27 mAP in default full setting, 2.02 mAP in default rare setting, 2.55 mAP in default non-rare setting, 2.42 mAP in known object full setting, 1.11 mAP in known rare settings and 2.79 mAP in known object non-rare setting. In addition, FGAHOI with Swin-Large*+ backbone exceeds an especially significant gain of 1.75 mAP in default full

TABLE 5: Comparison on ten intervals of the two proposed challenges. We divide the HICO-DET dataset into ten intervals based on each of the two challenges and compare the performance of QAHOI and FGAHOI on each interval. Kindly refer to Sec. 5.5 for more details.

| Challenge | Method | Backbone | $mAP_{role}$ (↑) | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $IMI_0$ | $IMI_1$ | $IMI_2$ | $IMI_3$ | $IMI_4$ | $IMI_5$ | $IMI_6$ | $IMI_7$ | $IMI_8$ | $IMI_9$ |
| UDA | QAHOI | Swin-Tiny | 16.35 | 24.72 | 29.24 | 34.79 | 38.70 | 46.21 | 53.13 | 47.60 | 58.66 | 60.19 |
| | | Swin-Large$_+^*$ | 20.53 | 33.58 | 41.11 | 45.41 | 45.44 | 56.43 | 56.25 | 63.53 | 71.12 | 75.08 |
| | FGAHOI | Swin-Tiny | 19.74 | 29.85 | 32.20 | 39.46 | 40.54 | 48.55 | 51.32 | 46.50 | 66.44 | 78.17 |
| | | Swin-Large$_+^*$ | 23.69 | 35.85 | 42.51 | 50.50 | 46.89 | 56.95 | 56.33 | 63.04 | 75.70 | 79.42 |
| LDVM | QAHOI | Swin-Tiny | 1.33 | 4.43 | 2.57 | 5.00 | 8.06 | 17.87 | 22.81 | 29.25 | 34.03 | 42.29 |
| | | Swin-Large$_+^*$ | 0.82 | 4.08 | 2.56 | 7.53 | 11.42 | 22.87 | 30.94 | 41.38 | 45.31 | 60.15 |
| | FGAHOI | Swin-Tiny | 2.50 | 4.15 | 3.34 | 7.58 | 9.83 | 21.61 | 27.64 | 33.07 | 38.31 | 45.07 |
| | | Swin-Large$_+^*$ | 1.44 | 4.32 | 4.57 | 7.81 | 11.82 | 24.92 | 32.50 | 43.66 | 47.26 | 60.55 |

TABLE 6: We carefully ablate each of the constituent component of FGAHOI. The middle results denote the role mAP. The results in the top right corner represent the performance improvement compared to QAHOI. The results in the bottom right corner represent the performance improvement compared to the baseline. Kindly refer to Sec. 5.7.1 for more details.

| Method | Merging Mechanism | | | Default | | | Known Object | | |
|---|---|---|---|---|---|---|---|---|---|
| | Hierarchical Spatial-Aware | Task-Aware | Full ↑ | Rare ↑ | Non-Rare ↑ | Full ↑ | Rare ↑ | Non-Rare ↑ |
| QAHOI | - | | - | 27.67 | 20.22 | 29.69 | 30.06 | 22.95 | 32.18 |
| FGAHOI | ✗ | ✗ | $28.45_{(\ -\ )}^{(+0.78)}$ | $21.07_{(\ -\ )}^{(+0.85)}$ | $30.66_{(\ -\ )}^{(+0.97)}$ | $31.08_{(\ -\ )}^{(+1.02)}$ | $24.02_{(\ -\ )}^{(+1.01)}$ | $33.19_{(\ -\ )}^{(+1.07)}$ |
| | ✔ | ✗ | $29.60_{(+1.15)}^{(+1.93)}$ | $\mathbf{22.39}_{(+1.32)}^{(+2.17)}$ | $31.76_{(+1.10)}^{(+2.07)}$ | $32.07_{(+0.99)}^{(+2.01)}$ | $\mathbf{24.48}_{(+0.46)}^{(+1.53)}$ | $34.34_{(+1.15)}^{(+2.16)}$ |
| | ✗ | ✔ | $29.32_{(+0.87)}^{(+1.65)}$ | $22.34_{(+1.27)}^{(+2.12)}$ | $31.41_{(+0.75)}^{(+1.72)}$ | $31.81_{(+0.73)}^{(+1.75)}$ | $24.30_{(+0.28)}^{(+1.35)}$ | $34.05_{(+0.86)}^{(+1.87)}$ |
| | ✔ | ✔ | $\mathbf{29.94}_{(+1.49)}^{(+2.27)}$ | $22.24_{(+1.17)}^{(+2.02)}$ | $\mathbf{32.24}_{(+1.58)}^{(+2.55)}$ | $\mathbf{32.48}_{(+1.40)}^{(+2.42)}$ | $24.16_{(+0.14)}^{(+1.21)}$ | $\mathbf{34.97}_{(+1.78)}^{(+2.79)}$ |

setting, 1.49 mAP in default rare setting, 1.82 mAP in default non-rare setting, 1.7 mAP in known object full setting, 0.92 mAP in known rare settings and 1.93 mAP in known object non-rare setting.

### 5.4.2 HOI-SDC

On the dataset we propose, *i.e.*, HOI-SDC, we compare FGAHOI with QAHOI and ablate each component of FGA-HOI (As shown in Table.3). The backbone is set to Swin-Tiny. The baseline exceeds QAHOI an especially significant gain of 1.63 mAP. HSAM and TAM improve a significant gain of 0.73 and 0.66 mAP, respectively. Benefit from the MSS, HSAM and TAM, FGAHOI achieve 22.25 mAP on HOI-SDC.

### 5.4.3 V-COCO

We compare FGAHOI with the state-of-the-art methods on V-COCO dataset and report the results in Table.4. In comparison to QAHOI, FGAHOI only exceeds a small margin. This phenomenon is mainly caused by too little training data in the dataset. We investigate that FGAHOI cannot adequately perform when the training data is not sufficient due to the complex task requirements. In addition, we investigate the transformer backbone is still superior to CNN backbone in this case.

### 5.5 Sensitivity Analysis for UDA and LDVM

According to the two proposed challenges, we divide the HICO-DET into ten intervals. At each intervals, we compare FGAHOI and QAHOI with Swin-Tiny, Large$_+^*$ backbone, respectively (As shown in Table.5). When compared between each interval of UDA and LDVM, we investigate that the difficulty of HOI detection decreases as the interval level increases. This justifies the original design. Thus, it is imperative to consider ability of the model to address these two challenges when proposing novel frameworks for HOI detection. In the comparison between FGAHOI and QAHOI, the results demonstrate that FGAHOI has better capability for uneven distributed area and long distance visual modeling of human-object pairs.

### 5.6 Qualitative Analysis

#### 5.6.1 Visualized Results

In order to demonstrate our model, several representative HOI predictions are visualized. As shown in Fig.5, our model can pinpoint HOI instances from noisy backgrounds and excels at detecting various complicated HOIs, including one object interacting with different humans, one human engaging in multiple interactions with various objects, multiple interactions within a single pair, and multiple humans engaging in various interactions with various objects. In addition, our model is good at long-range visual modelling, withstanding the impacts of hostile environments and small target identification. Fig.6 (a) illustrates that FGAHOI has excellent long-range visual modelling capabilities and can accurately identify interactions between human-object pairs far from each other. As Fig.6 (b) shows, our model has
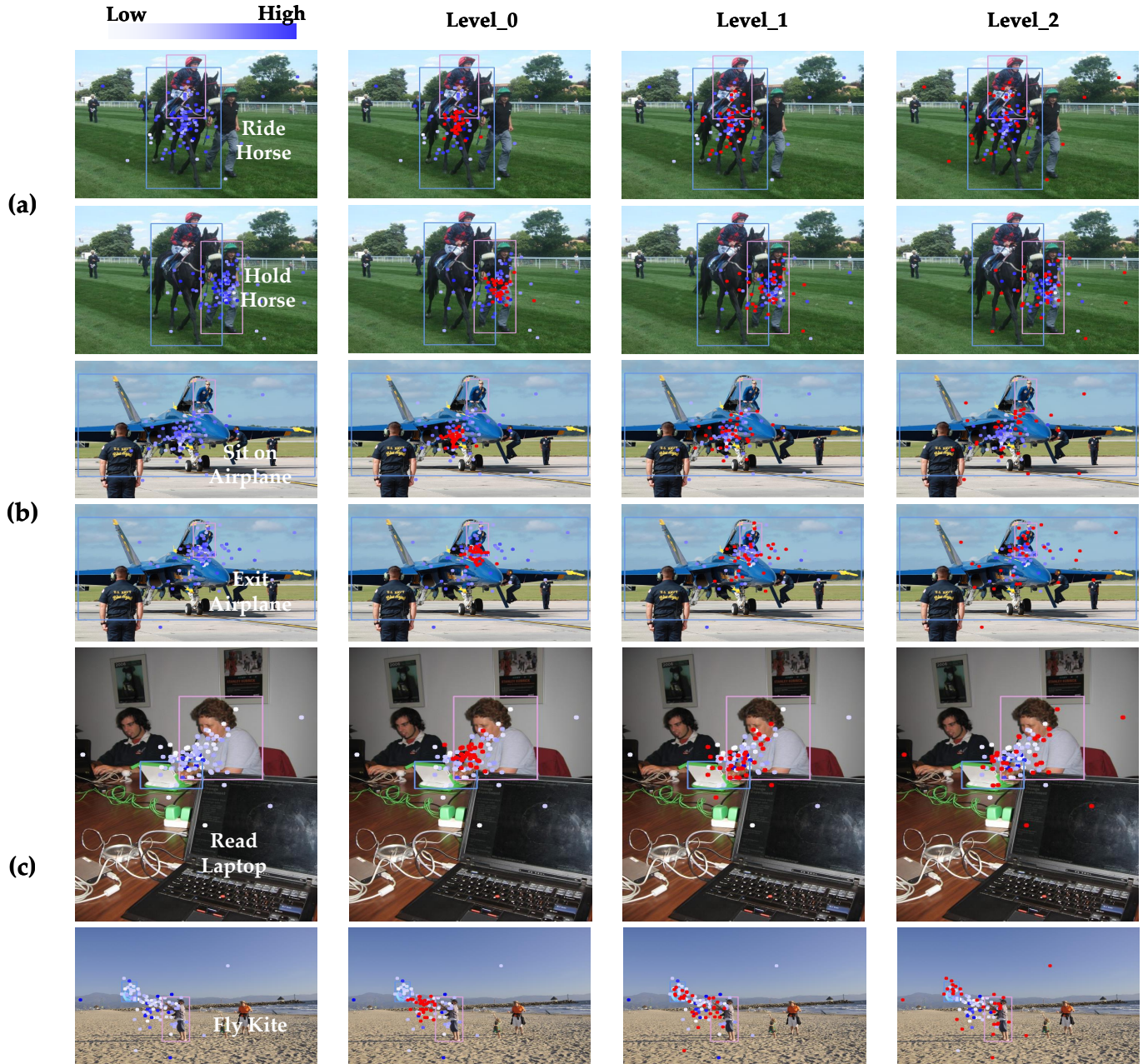
Fig. 10: Visualization of fine-grained anchors in the decoding phase, Level_0, Level_1 and Level_2 represent the features at different scales respectively, the color of the blue dots from light to dark represents the degrees of attention of the fine-grained anchors and red dots represent the positions of interest of fine-grained anchors in current scale features. Kindly refer to Sec. 5.6.2 for more details.

outstanding robustness and can effectively resist disruption from harsh environmental factors, including blurring, blocking and glare. Fig.6 (c) demonstrates the superior capabilities of FGAHOI to identify small HOI instances.

### 5.6.2   What do the fine-grained anchors look at?

As shown in Fig.7, we compare the fine-grained anchors of FGAHOI and QAHOI. First two HOI instances (*i.e*, hold sport ball and ride motorcycles) exhibit that FGAHOI could better focus on humans, objects and the interaction areas rather than noisy backgrounds. The fourth head of FGAHOI still focuses on the HOI instance, while QAHOI focuses

on the background. When detecting instance with a long distance between human and object, FGAHOI could focus on the right position, while QAHOI is like a chicken with its head cut off (As shown in the last HOI instance).

To exhibit the effectiveness of the fine-grained anchors for identifying HOI instances and demonstrate the working mechanism of fine-grained anchors, we visualize the fine-grained anchors of the feature maps at different scales in the decoding phase. In Fig.10 (a), we visualize the instances of two different humans and one object. As shown in Fig.10 (b), even for exactly the same human-object pair, the areas of focus vary from one interaction to another. In Fig.10 (c),

Fig. 11: Visualization of several representative interactive actions and the corresponding fine-grained anchors. We only visualize a single representative interactive action for one human-object pair. Kindly refer to Sec. 5.6.2 for more details.

we show two instances contain short and long distance between humans and objects, respectively. We investigate that the fine-grained anchors of low level feature map focus on small and fine-grained areas. They play a major role in detecting close range and small HOI instances. The fine-grained anchors of high level feature maps focus on large and coarse-grained areas. It is necessary for detecting long distance and large HOI instances.

In order to explore what the fine-grained anchors focus on, we visualize several representative actions in Fig.11. Visualization shows that fine-grained anchors could concentrate attention precisely on the location where the interactive action is generated. For example, the fine-grained anchors mainly focus on the hand for 'text_on cell_phone', the mouth for 'eat orange' and the ear and the mouth for 'talk_on cell_phone'. For 'kick sports_ball', 'jump skateboard' and 'hop_on elephant', central areas of interest are around legs and feet, while fine-grained anchors primarily focuses on hands for 'carry handbag', 'repair hair_drier', 'hold cup', 'hold hotdog' and 'cut with kinfe'.

### 5.7 Ablation Study

In this subsection, a set of experiments are designed to clearly understand the contribution of each of the constituent components of the proposed methodology: **Merging mechanism**, **Multi-Scale Sampling Strategy** and **Stage-wise Training Strategy**. We conducted all experiments on the HICO-DET dataset.

#### 5.7.1 Ablating FGAHOI Components

To study the contribution of each of the merging mechanisms in FGAHOI, we design careful ablation experiments in Table.6. To ensure a fair comparison, the sampling sizes are all set to [1, 3, 5]. For the baseline which does not leverages the hierarchical spatial-aware and task-aware merging mechanism, we use the average and direct summation operation to merge the sampled features and connect embeddings. For the results in the table, the middle results denote the role mAP, the results in the top right corner represent the performance improvement compared to QAHOI and the results in the bottom right corner represent the performance improvement compared to the baseline. In comparison to row 1 (QAHOI), row 2 adds the multi-scale sampling strategy. The results demonstrate that adding the sampling strategy improves the ability of the model to detect HOI instances. The row 3 and 4 show that both hierarchical spatial-aware and task-aware merging mechanism make an essential contribution to the success of FGAHOI. The hierarchical spatial-aware merging mechanism, combined with the task-aware merging mechanism performs better together (row 5) than using either of them separately (row 3 and 4). Thus, each component in FGAHOI has a critical role to play in HOI detection.

#### 5.7.2 Sensitivity Analysis On Multi-Scale Sampling Sizes

Our multi-scale sampling strategy samples multi-scale features according to the pre-determined sampling sizes. We vary different sampling sizes to conduct the sensitivity

analysis for the sampling strategy and report the results in Table.7. We find that the sampling strategy is relatively stable. Changes in sampling sizes do not have a significant impact on the performance of FGAHOI. However, there is still a slight degradation in the performance of FGAHOI as the sample size increases. We investigate that as the sample size increases, too many background features around the fine-grained anchors are sampled, resulting in contamination of the sampled features and thus the performance of the model suffers. Hence, for validation, we set the sampling sizes to [1, 3, 5] in all our experiments, which is a sweet spot that balances performance.

TABLE 7: Comparison between different sampling sizes.

| Smpling Size | Default | | | Known Object | | |
|---|---|---|---|---|---|---|
| | Full | Rare | Non-Rare | Full | Rare | Non-Rare |
| [ 1, 3, 5 ] | **29.94** | 22.24 | **32.24** | 32.48 | 24.16 | **34.97** |
| [ 3, 5, 7 ] | 29.72 | **23.03** | 31.72 | 32.33 | **25.67** | 34.30 |
| [ 5, 7, 9 ] | 29.65 | 22.64 | 31.74 | 32.55 | 25.64 | 34.62 |

### 5.7.3 Training Strategies

As shown in Table.8, we leverage the stage-wise and end-to-end training strategy to train FGAHOI, respectively. In the end-to-end training strategy, we train FGAHOI for 150 epochs and the learning rate drop is carried out at the 120th epoch. The stage-wise training strategy promotes 5.96 mAP for default full setting, 4.61 for default rare, 6.36 for default non-rare, 6.04 for known object full, 4.65 for known object rare and 6.46 mAP for known object non-rare setting. In comparison to the end-to-end training strategy, we investigate that the stage-wise training strategy reduces the learning difficulty of the FGAHOI and clarify the learning direction of the model by emphasizing it to learn what it needs at each stage.

TABLE 8: Comparison between Stage-Wise and End-to-End training approach.

| Training Strategy | Default | | | Known Object | | |
|---|---|---|---|---|---|---|
| | Full | Rare | Non-Rare | Full | Rare | Non-Rare |
| Stage-Wise | **29.94** | **22.24** | **32.24** | **32.48** | **24.16** | **34.97** |
| End-to-End | 23.98 | 17.63 | 25.88 | 26.44 | 19.51 | 28.51 |

## 6 CONCLUSION

In this paper, we propose a novel transformer-based human-object interaction detector (FGAHOI) which leverages the input features to generate fine-grained anchors for protecting the detection of HOI instances from noisy backgrounds. We propose a novel training strategy where each component of the model is trained sequentially to clarify the training direction at each stage, for maximizing the savings of the training cost. We propose two novel metrics and a novel dataset, *i.e.*, HOI-SDC for the two challenges (Uneven Distributed Area in Human-Object Pairs and Long Distance Visual Modeling of Human-Object Pairs) of detecting HOI instances. Our extensive experiments on three benchmarks: HICO-DET, HOI-SDC and V-COCO, demonstrate the effectiveness of the proposed FGAHOI. Specifically, FGAHOI outperforms all existing state-of-the-art methods by a large margin.

## REFERENCES

[1] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448, 2015.
[2] Z. Li and F. Zhou, "Fssd: feature fusion single shot multibox detector," *arXiv preprint arXiv:1712.00960*, 2017.
[3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
[4] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 580–587, 2014.
[5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779–788, 2016.
[6] M. Chen, Y. Liao, S. Liu, Z. Chen, F. Wang, and C. Qian, "Reformulating hoi detection as adaptive set prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9004–9013, 2021.
[7] Z. Hou, X. Peng, Y. Qiao, and D. Tao, "Visual compositional learning for human-object interaction detection," in *European Conference on Computer Vision*, pp. 584–600, Springer, 2020.
[8] P. Zhou and M. Chi, "Relation parsing neural network for human-object interaction detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 843–851, 2019.
[9] B. Kim, J. Lee, J. Kang, E.-S. Kim, and H. 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*, pp. 74–83, 2021.
[10] X. Zhong, C. Ding, X. Qu, and D. Tao, "Polysemy deciphering network for robust human–object interaction detection," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1910–1929, 2021.
[11] X. Zhong, X. Qu, C. Ding, and D. 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*, pp. 13234–13243, 2021.
[12] C. Gao, J. Xu, Y. Zou, and J.-B. Huang, "Drg: Dual relation graph for human-object interaction detection," in *European Conference on Computer Vision*, pp. 696–712, Springer, 2020.
[13] O. Ulutan, A. Iftekhar, and B. 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*, pp. 13617–13626, 2020.
[14] C. Gao, Y. Zou, and J.-B. Huang, "ican: Instance-centric attention network for human-object interaction detection," *arXiv preprint arXiv:1808.10437*, 2018.
[15] T. Wang, T. Yang, M. Danelljan, F. S. Khan, X. Zhang, and J. Sun, "Learning human-object interaction detection using interaction points," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4116–4125, 2020.
[16] A. Zhang, Y. Liao, S. Liu, M. Lu, Y. Wang, C. Gao, and X. Li, "Mining the benefits of two-stage and one-stage hoi detection," *Advances in Neural Information Processing Systems*, vol. 34, pp. 17209–17220, 2021.
[17] C. Zou, B. Wang, Y. Hu, J. Liu, Q. Wu, Y. Zhao, B. Li, C. Zhang, C. Zhang, Y. 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*, pp. 11825–11834, 2021.

[18] F. Z. Zhang, D. Campbell, and S. Gould, "Spatially conditioned graphs for detecting human-object interactions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 13319–13327, 2021.

[19] M. Tamura, H. Ohashi, and T. 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*, pp. 10410–10419, 2021.

[20] Y.-L. Li, S. Zhou, X. Huang, L. Xu, Z. Ma, H.-S. Fang, Y. Wang, and C. Lu, "Transferable interactiveness knowledge for human-object interaction detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3585–3594, 2019.

[21] G. Gkioxari, R. Girshick, P. Dollár, and K. He, "Detecting and recognizing human-object interactions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8359–8367, 2018.

[22] Y.-W. Chao, Y. Liu, X. Liu, H. Zeng, and J. Deng, "Learning to detect human-object interactions," in *2018 ieee winter conference on applications of computer vision (wacv)*, pp. 381–389, IEEE, 2018.

[23] T. Gupta, A. Schwing, and D. Hoiem, "No-frills human-object interaction detection: Factorization, layout encodings, and training techniques," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9677–9685, 2019.

[24] B. Wan, D. Zhou, Y. Liu, R. Li, and X. He, "Pose-aware multi-level feature network for human object interaction detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9469–9478, 2019.

[25] A. Bansal, S. S. Rambhatla, A. Shrivastava, and R. Chellappa, "Detecting human-object interactions via functional generalization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 10460–10469, 2020.

[26] H.-S. Fang, Y. Xu, W. Wang, X. Liu, and S.-C. Zhu, "Learning pose grammar to encode human body configuration for 3d pose estimation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018.

[27] H.-S. Fang, G. Lu, X. Fang, J. Xie, Y.-W. Tai, and C. Lu, "Weakly and semi supervised human body part parsing via pose-guided knowledge transfer," *arXiv preprint arXiv:1805.04310*, 2018.

[28] Y. Xiu, J. Li, H. Wang, Y. Fang, and C. Lu, "Pose flow: Efficient online pose tracking," *arXiv preprint arXiv:1802.00977*, 2018.

[29] S. Qi, W. Wang, B. Jia, J. Shen, and S.-C. Zhu, "Learning human-object interactions by graph parsing neural networks," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 401–417, 2018.

[30] B. Kim, T. Choi, J. Kang, and H. J. Kim, "Uniondet: Union-level detector towards real-time human-object interaction detection," in *European Conference on Computer Vision*, pp. 498–514, Springer, 2020.

[31] Y. Liao, S. Liu, F. Wang, Y. Chen, C. Qian, and J. 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*, pp. 482–490, 2020.

[32] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*, pp. 213–229, Springer, 2020.

[33] J. Chen and K. Yanai, "Qahoi: Query-based anchors for human-object interaction detection," *arXiv preprint arXiv:2112.08647*, 2021.

[34] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.

[35] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," *Advances in neural information processing systems*, vol. 19, 2006.

[36] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.

[37] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," *arXiv preprint arXiv:1910.09217*, 2019.

[38] Y.-W. Chao, Y. Liu, X. Liu, H. Zeng, and J. Deng, "Learning to detect human-object interactions," in *2018 ieee winter conference on applications of computer vision (wacv)*, pp. 381–389, IEEE, 2018.

[39] S. Gupta and J. Malik, "Visual semantic role labeling," *arXiv e-prints*, 2015.

[40] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.

[41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[42] Y. LeCun, Y. Bengio, *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.

[43] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021.

[44] X. Chen, F. Wei, G. Zeng, and J. Wang, "Conditional detr v2: Efficient detection transformer with box queries," *arXiv preprint arXiv:2207.08914*, 2022.

[45] Y. Wang, X. Zhang, T. Yang, and J. Sun, "Anchor detr: Query design for transformer-based detector," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 36, pp. 2567–2575, 2022.

[46] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "Dab-detr: Dynamic anchor boxes are better queries for detr," *arXiv preprint arXiv:2201.12329*, 2022.

[47] G. Zhang, Z. Luo, Y. Yu, K. Cui, and S. Lu, "Accelerating detr convergence via semantic-aligned matching," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 949–958, 2022.

[48] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

[49] F. Iandola, M. Moskewicz, S. Karayev, R. Girshick, T. Darrell, and K. Keutzer, "Densenet: Implementing efficient convnet descriptor pyramids," *arXiv preprint arXiv:1404.1869*, 2014.

[50] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European conference on computer vision*, pp. 483–499, Springer, 2016.

[51] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, "Cswin transformer: A general vision transformer backbone with cross-shaped windows," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12124–12134, 2022.

[52] J. Yang, C. Li, P. Zhang, X. Dai, B. Xiao, L. Yuan, and J. Gao, "Focal self-attention for local-global interactions in vision transformers," *arXiv preprint arXiv:2107.00641*, 2021.

[53] Z. Xia, X. Pan, S. Song, L. E. Li, and G. Huang, "Vision transformer with deformable attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4794–4803, 2022.

[54] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22–31, 2021.

[55] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[56] C.-F. Chen, R. Panda, and Q. Fan, "Regionvit: Regional-to-local attention for vision transformers," *arXiv preprint arXiv:2106.02689*, 2021.

[57] Y. Wang, R. Huang, S. Song, Z. Huang, and G. Huang, "Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition," *Advances in Neural Information Processing Systems*, vol. 34, pp. 11960–11973, 2021.

[58] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.

[59] X. Dai, Y. Chen, B. Xiao, D. Chen, M. Liu, L. Yuan, and L. Zhang, "Dynamic head: Unifying object detection heads with attentions," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7373–7382, 2021.

[60] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 734–750, 2018.

[61] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.

[62] Y.-L. Li, L. Xu, X. Liu, X. Huang, Y. Xu, S. Wang, H.-S. Fang, Z. Ma, M. Chen, and C. Lu, "Pastanet: Toward human activity knowledge engine," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 382–391, 2020.

[63] Y.-W. Chao, Z. Wang, Y. He, J. Wang, and J. Deng, "Hico: A benchmark for recognizing human-object interactions in images," in *Proceedings of the IEEE international conference on computer vision*, pp. 1017–1025, 2015.

[64] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.