

Decoupling Object Detection from Human-Object Interaction Recognition

Ying Jin^{†‡} Yinpeng Chen[†] Lijuan Wang[†] Jianfeng Wang[†]
 Pei Yu[†] Lin Liang[†] Jenq-Neng Hwang[‡] Zicheng Liu[†]

Microsoft[†], University of Washington[‡]

{ying.jin, yiche, lijuanw, jianfw, peyu, lliang, zliu}@microsoft.com

{jinying, hwang}@uw.edu

Abstract

We propose *DEFR*, a *DE*tection-*FR*ee method to recognize Human-Object Interactions (HOI) at image level **without** using object location or human pose. This is challenging as the detector is an integral part of existing methods. In this paper, we propose two findings to boost the performance of the detection-free approach, which significantly outperforms the detection-assisted state of the arts. Firstly, we find it crucial to effectively leverage the semantic correlations among HOI classes. Remarkable gain can be achieved by using language embeddings of HOI labels to initialize the linear classifier, which encodes the structure of HOIs to guide training. Further, we propose *Log-Sum-Exp Sign (LSE-Sign)* loss to facilitate multi-label learning on a long-tailed dataset by balancing gradients over all classes in a softmax format. Our detection-free approach achieves 65.6 mAP in HOI classification on HICO, outperforming the detection-assisted state of the art (SOTA) by 18.5 mAP, and 52.7 mAP in one-shot classes, surpassing the SOTA by 27.3 mAP. Different from previous work, our classification model (*DEFR*) can be directly used in HOI detection without any additional training, by connecting to an off-the-shelf object detector whose bounding box output is converted to binary masks for *DEFR*. Surprisingly, such a simple connection of two decoupled models achieves SOTA performance (32.35 mAP).

1. Introduction

Human-Object Interaction (HOI) recognition has drawn significant interest for its important role in scene understanding. HOI recognition aims to retrieve multiple $\langle verb, object \rangle$ pairs in the entire image. Previous studies on image-level HOI classification [5, 9, 12, 14, 29, 36, 45] mostly rely on object detectors to detect humans and objects in the first stage, and then infer interactions between human-object pairs. The pipelines are becoming more complicated as

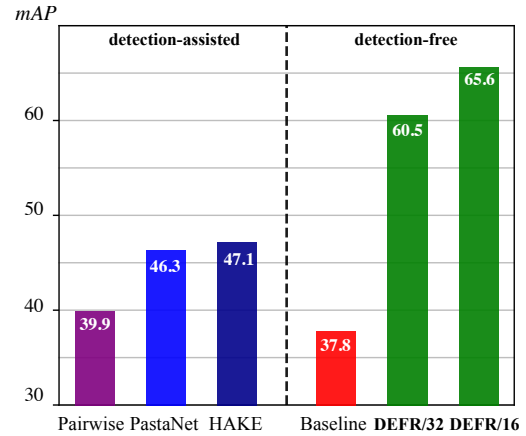


Figure 1. **Detection-Assisted v.s Detection-Free** HOI classifiers. Although the detection-free baseline (using ImageNet pre-trained ViT-B/32 [8] as backbone) degrades severely from the detection-assisted HAKE [28], DEFR, our method, achieves significant improvement by leveraging language embedding as classifier initialization and the proposed LSE-Sign loss. DEFR/16 (ViT-B/16 backbone) achieves 65.6 mAP on the HICO dataset.

human pose and body part states [29, 35] are added as additional features. In contrast, this paper shifts focus to a *detection-free* solution without affecting the performance.

By eliminating the use of object and human keypoint detections, our detection-free approach has a significantly simplified pipeline. However, the goal is nontrivial as the positions of the HOI-related human and object are unknown. A naive baseline with a vision transformer backbone pre-trained on ImageNet-1K experiences a considerable performance degradation compared to the state of the art [28] (37.8 mAP vs. 47.1 mAP in Figure 1).

In this paper, we show that the performance of a detection-free approach can be significantly boosted by using *language embedding as classifier initialization* and the proposed *log-sum-exp sign* loss function. The former not only encodes the structure of HOI classes to guide training, but also provides

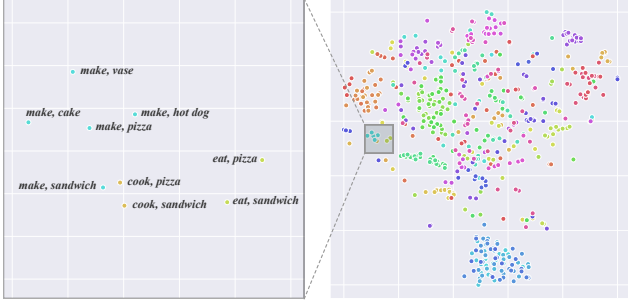


Figure 2. **TSNE visualization of BERT-encoded HOI text.** Each point represents the language embedding of an HOI class. **Left:** zoomed view of a cluster of HOIs. HOIs with closer semantic meanings (e.g. $\langle \text{make}, \text{hot dog} \rangle$, $\langle \text{make}, \text{pizza} \rangle$, $\langle \text{cook}, \text{pizza} \rangle$) locate closely. **Right:** 600 HOI classes from the HICO dataset sorted by the verb. The figure shows a clear clustering effect, which captures the structure of HOI classes.

proper weight initialization for few-shot classes in the linear classifier. The latter successfully handles multiple HOI labels per image in a long-tailed dataset. Our method, named DETection FRee (DEFR) HOI recognition, is based upon these two pillars.

We find using language embeddings as classifier initialization is crucial for HOI recognition. This is because HOI classes are semantically highly correlated, since HOIs are combined phrases of verbs and objects. This distinguishes HOI from traditional image classification tasks like ImageNet. Language models like BERT can encode the semantic structure of HOI classes effectively (Fig. 2). To leverage language, we use language embeddings of HOI labels as weight initialization in the classification head, which not only guides the training but also provides proper weight initialization for few-shot classes. We successfully boost the performance from 37.8 mAP to 50.6 mAP with ViT-B/32 backbone pre-trained on ImageNet-1K, and the gain is maximized when the language model is jointly trained with the image encoder. This technique, later referred to as *Language Embedding Initialization*, is generally effective for different backbones and pre-trainings.

In terms of the loss function, we propose Log-Sum-Exp Sign (LSE-Sign) loss to handle multiple labels per image. Compared to binary entropy loss which considers each class independently, LSE-Sign loss enables interaction between classes by balancing the gradients over all classes in a softmax format. This encourages more attention towards the class with maximum loss, and benefits the learning on a long-tailed dataset.

Based upon these two new findings, our detection-free HOI classification model outperforms previous detection-assisted approaches [9, 28, 29] by a clear margin. Our best model achieves 65.6 mAP on the HICO [5] dataset, surpassing the state of the art (SOTA) [28] that uses both object

and human keypoint detections, by 18.5 mAP. On few-shot subsets, our model achieves 52.7 mAP for 1-shot and 56.9 mAP for 5-shot, respectively 27.3 and 24.4 mAP gains over the existing SOTA. Furthermore, we tackle HOI detection by connecting our classification model with an off-the-shelf object detector to recognize the regional HOI. Surprisingly, the model **without** additional training achieves SOTA performance (32.35 mAP) on the HICO-DET dataset, outperforming the SOTA [49] that is trained on HICO-DET.

2. Related work

HOI Classification: For HOI classification, a key challenge is the co-occurrence of multiple humans and objects when the positions of HOIs in an image are not labeled. Existing work [9, 14, 29, 36] depend on object or human keypoint detectors and use Multiple Instance Learning (MIL) to train on non-localized HOIs [37]. Girdhar et al. [12] utilizes attention maps guided by human pose and removes MIL. PastaNet [9] achieve state-of-the-art performance by conducting body part-level action classification to parse the HOI.

HOI Detection: Instance level HOI detection on HICO-DET [4] is attracting growing interest. Existing work can be categorised into three streams. Two-stage methods [10, 11, 13, 22, 29, 33, 50] require object detection prior to HOI classification to extract regional features. Graph neural networks are often used [50] to classify the verbs between human-object pairs. For training, two-stage methods usually initialize their backbone from a pre-trained object detector. One-stage methods [20, 30, 51] execute object detection and HOI detection in parallel and match them afterwards. Recent studies [6, 21, 42, 49, 52] achieve end-to-end HOI detection based on DETR [3] and benefit from the wider perception field of transformers [42]. Language is used as priors or features. [22] considers verb-object co-occurrence priors and [10, 28] use text embeddings of detected objects or part-level body actions as features.

Scene Graph Generation: As in connection with HOI, scene graphs [25] are also pairwise relationships. A scene graph is a graph representation in which the nodes are objects and edges are relationships. The IMP [45] method detects objects first and then region features are iteratively fine-tuned in a graphical network. To reduce the compute, [26] constructs sparsely connected sub-graphs, and [46] applies a Relationship Proposal Network. [48] reveals the strong statistical bias in the Visual Genome [25] dataset, and [43] suggests unbiased post-processing to neutralize the model from biased training. [7] and [47] use frequency priors from the dataset and text as external knowledge to improve performance.

3. Detection-free HOI recognition

In this section, we introduce our detection-free solution (DEFR) for Human-Object Interaction (HOI) recognition.

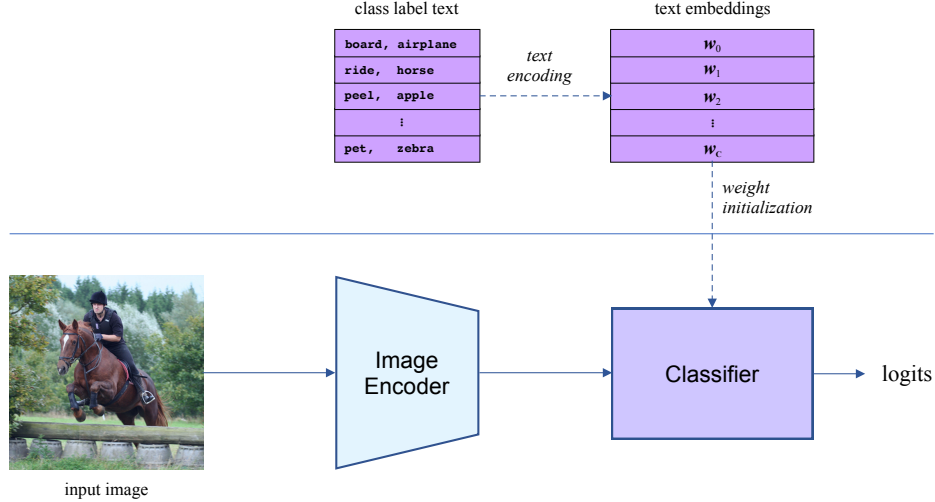


Figure 3. **DEtection FRee (DEFR) HOI recognition pipeline.** It has an image encoder and a linear classifier. The weight of the linear classifier is initialized by w , the text embeddings of HOI classes encoded by a language model. We call this *language embedding initialization*, detailed in [subsection 3.2](#). Compared to detection-assisted approaches, DEFR significantly simplifies the pipeline.

The goal of DEFR is to eliminate the dependency of object and human keypoint detection without performance degradation. Surprisingly, we find that performance can be boosted by using language embedding of HOI classes as classifier initialization and the proposed log-sum-exp sign (LSE-Sign) loss function.

3.1. Detection-free HOI recognition pipeline

Our DEFR method has a simple pipeline (shown in [Figure 3](#)), including a vision transformer [8] as the backbone and a linear layer as the classifier. The classifier is initialized with language embeddings, instead of the conventional random initialization methods [15, 16]. Compared to detection-assisted methods [9, 10, 28, 29], our approach simplifies the pipeline significantly and enables learning in an end-to-end manner.

With the simple pipeline, the problem becomes challenging due to the elimination of detection input, which provides accurate spatial information about humans and objects. To validate this, we fine-tune a vision transformer (ViT-B/32) backbone pre-trained on ImageNet-1K on the HICO [5] dataset. A severe degradation (to 37.8 mAP) is observed when compared with the detection-assisted approaches (47.1 mAP in HAKE [28]).

3.2. Language Embedding Initialization

HOI classes have strong semantic correlations as they are combinatorial phrases of verbs and objects. Such correlation differentiates HOI recognition from traditional classification tasks like ImageNet classification, and is vital to the performance of HOI recognition. Language models can effectively encode the structure of HOI classes, which can not be effectively learned by the image backbone along (see [Fig. 2](#)).

Based on this observation, we use language embedding to initialize the linear classifier layer to guide training and successfully boost the model performance. We refer to this technique as *Language Embedding initialization*.

Specifically, we convert the HOI classes (e.g. “<ride, bicycle>”) to prompts (e.g. “a person riding a bicycle”) and generate language embeddings with a text encoder. The embeddings are normalized and used as the initial weight in the final linear layer which produces class logits. During training, the output logit for the i^{th} class is the dot product of the image feature and w_i , a row vector in the classifier’s weight w , if bias is not used (see [Fig. 3](#)). Since dot product is the unnormalized cosine similarity, w_i is often considered as the proxy for a given class [41, 44]. Using language embeddings of HOI classes to initialize w (proxies) in the classifier not only encodes the structure of HOI classes to guide training, but also provides proper weight initialization for few-shot classes.

In our scenario, language embedding initialization using BERT fires a performance boost of 12.8 mAP for an ImageNet-1K pre-trained ViT-B/32 backbone. More importantly, the applicability generalizes to backbone architectures (ViT and ResNet [17]) pre-trained in different ways. The gain is maximized when the language model is jointly trained with the image encoder (e.g. CLIP). [Figure 4](#) visualizes the linear classifier’s weight vectors per class before (left-column) and after fine-tuning with both CLIP (middle-column) and ImageNet-1K (right-column) pre-trained backbones. The weight vectors initialized with language embeddings maintains clustered after fine-tuning. However, this structure is difficult to learn when fine-tuned with random initialization, together with lower performance (shown in [Figure 4](#) bottom row).



Figure 4. **TSNE visualization of classifier weights.** **Left Column:** the CLIP-encoded text embeddings of 600 HOI class labels. Each point represents the embedding of an HOI class. Text embeddings encode the structure of HOI classes: as is shown in the zoomed view, points that represent similar semantic meanings locate closely. **Middle Column:** the weight vectors of 600 HOI classes (i.e. row vectors in linear classifier’s weight matrix) after fine-tuning, when CLIP pre-trained ViT-B/32 backbone is used. The top and bottom use language embedding initialization and random initialization, respectively. **Right Column:** the weight vectors of 600 HOI classes after fine-tuning, when ImageNet-1K pre-trained ViT-B/32 backbone is used. The top and bottom use text embedding and random initialization for the linear classifier, respectively. Text embeddings are clearly clustered before fine-tuning, and the overall structure changes slightly after fine-tuning. However, this clustered structure is not clearly learned when fine-tuned with random initialization, reflected in lower model performance (see mAP results in the plots).

3.3. Log-Sum-Exp Sign (LSE-Sign) loss

We propose a log-sum-exp sign (LSE-Sign) loss function to facilitate multi-class learning, as an image usually contains multiple interactions (e.g. *cut carrot*, *hold carrot*, *peel carrot*). Let \mathbf{x} denote a feature vector from the backbone and $\mathbf{y} = \{y_1, y_2, \dots, y_C\}$ denote multiple class labels, where C is the number of HOI classes, and $y_i \in \{1, -1\}$, indicating positive and negative classes, respectively. We first normalize the logit output per class and re-scale it as follows:

$$s_i = \gamma \frac{\mathbf{x}^T \mathbf{w}_i}{\|\mathbf{x}\| \|\mathbf{w}_i\|}, \quad (1)$$

where \mathbf{w}_i is the i^{th} row in the weight matrix of the linear classifier, corresponding to the i^{th} HOI class, and γ is a scalar hyper-parameter, controlling the output range. It is

equivalent to scaling cosine similarity such that the output value is between $\pm\gamma$. The losses of both positive and negative classes can be unified into $e^{-y_i s_i}$, where the label y_i controls the sign. The overall loss is defined as the log-sum-exp function of $-y_i s_i$ as follows:

$$\mathcal{L} = \log \left(1 + \sum_{i=1}^C e^{-y_i s_i} \right) \quad (2)$$

where the constant term 1 in the log function sets the zero lower bound for the loss. Log-sum-exp is a smooth approximation of the maximum function, and its gradient is the softmax function as follows:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial s_i} &= \frac{-y_i e^{-y_i s_i}}{1 + \sum_{j=1}^C e^{-y_j s_j}} \\ &= \frac{-1_{i \in pos} e^{-s_i} + 1_{i \in neg} e^{s_i}}{1 + \sum_{i \in pos} e^{-s_i} + \sum_{i \in neg} e^{s_i}},\end{aligned}\quad (3)$$

where $1_{i \in pos}$ is a Dirac delta function that returns 1 if the i^{th} class is positive and 0 otherwise. Compared to the binary cross entropy loss and focal loss that consider each class separately, LSE-Sign loss considers the dependency across classes as the magnitude of the gradients are normalized over all classes and distributed by a softmax function. This facilitates multi-label learning on a long-tailed dataset as it encourages learning of classes with larger loss values and suppresses the learning of classes with smaller loss values.

4. Experiments

We evaluate the proposed DEFR on both HOI classification and HOI detection.

4.1. HOI Classification

Dataset: For image-level HOI classification, we conduct experiments on two commonly used datasets: HICO and MPII. The HICO dataset [5] contains 600 HOI categories which have 117 unique verbs and 80 COCO [32] object classes. Each image may contain multiple HOI classes and multiple human-object pairs. The training set has 38,116 images and test set has 9,658 images. Following existing methods, we randomly reserve 10% images from the training set for validation and report performance on the test set. The MPII dataset [1] contains 15,205 training images and 5,708 test images. Unlike HICO, each image is labeled with only one of 393 interaction classes. We follow [9, 36] to report performance on the validation set that contains 6,987 images.

Pre-training: The backbone is ViT-B/32, and we investigate three different strategies to pre-train the backbone at resolution 224.

1. Image classification task on ImageNet-1K or ImageNet-21K referred to as CLS1K and CLS21K, respectively.
2. Masked language modeling task motivated by [23], where the network input includes both an image and the associated text. Both modalities are fully attended to each other in each transformer block. The pre-training datasets are Google Conceptual Captions [40], SBU [38], COCO [32] and Visual Genome [25]. This is referred to as MLM.
3. Image-text contrastive learning task based on CLIP [39]. The image encoder is jointly trained with a text encoder and the two modalities are contrasted on the encoder output representations. Only the image encoder is used as the backbone. Here, we directly use

Table 1. **Comparison with the state of the art on HICO.** Dependencies (additional input) required by existing methods include Bbox: object detection, Pose: human keypoints, PaSta [29]: additional training data of part level actions. We report the performance of DEFR/32 (ViT-B/32 backbone) and DEFR/16 (ViT-B/16 backbone). Compared to detection-assisted approaches, DEFR removes detection while achieving a significant boost in accuracy.

Method	Dependencies			mAP
	Bbox	Pose	PaSta	
R*CNN [14]	✓			28.5
Girdhar <i>et al.</i> [12]		✓		34.6
Mallya <i>et al.</i> [36]	✓			36.1
Pairwise-Part [9]	✓	✓		39.9
PastaNet [29]	✓	✓	✓	46.3
HAKI [28]	✓	✓	✓	47.1
DEFR/32				60.5
DEFR/16				65.6

the released CLIP model¹, and reuse the term CLIP to refer to the image encoder as pre-training.

Fine-tuning: All models have a ViT-B/32 backbone fine-tuned at resolution 672 with the AdamW [24] optimizer without weight decay. We use a batch size of 128 on 8 V100 GPUs. We set the base learning rate as 1.5e-5 for CLIP pre-trained backbones and 1e-4 otherwise, and use cosine scheduling with warm restarts [34] every 5 epochs. The best model is fine-tuned for 10 epochs. Data augmentation of random color jittering, horizontal flipping, and resized cropping is used. To reduce class imbalance, we adopt over-sampling so that each class has at least 40 samples per epoch.

4.2. Experiment results on HOI classification

HICO dataset: Table 1 compares our detection-free method (DEFR) with the prior works that need assistance from object detection or human keypoint detection. Our method achieves significantly better accuracy without detection assist. Specifically, DEFR achieves 65.6 mAP on the HICO, gaining over the state-of-the-art HAKI [28] by 18.5 mAP.

Few-shot analysis: Our model outperforms existing methods considerably in few-shot subsets on HICO, by leveraging language embedding initialization and LSE-Sign loss. As shown in Tab. 2, DEFR achieves 52.7 mAP in one-shot classes, gaining over the detection-assisted SOTA [28] by 27.3 mAP.

MPII dataset: we evaluate HOI classification on the MPII dataset in addition to HICO. As shown in Tab. 3, our model achieves SOTA performance of 55.3 mAP. Importantly, DEFR with ResNet101 backbone still outperforms [9] which is detection-assisted by a large margin. This proves that our

¹<https://github.com/openai/CLIP>

Table 2. **Few-shot performance** evaluated on HICO. Few@i means classes that the number of training images is i . The number of HOI classes for Few@1, 5, 10 are 49, 125 and 162, respectively.

Method	mAP	Few@1	Few@5	Few@10
Pairwise-Part [9]	39.9	13.0	19.8	22.3
PastaNet [29]	46.3	24.7	31.8	33.1
HAKE [28]	47.1	25.4	32.5	33.7
DEFR/16	65.6	52.7	56.9	57.2

Table 3. **Comparison on the MPII dataset.** We additionally apply the proposed approach on ResNet101 backbone for fair comparison. We follow [9, 14, 36] to report performance on the validation set that contains 6,987 images. Different from HICO, the MPII dataset has only one interaction label per image.

Method	Backbone	mAP
R*CNN [14]	VGG16	21.7
Girdhar <i>et al.</i> [12]	ResNet101	30.6
Pairwise-Part [9]	ResNet101	32.0
DEFR	ResNet101	43.6
DEFR	ViT-B/16	55.3

Table 4. **The path from baseline to DEFR** evaluated on HICO. The baseline uses ImageNet-1K pre-trained ViT-B/32 as the backbone, random initialization for the classifier, and binary cross entropy loss. Language embeddings are generated by BERT and CLIP’s text encoder as classifier’s weight initialization.

	Pre-training	LSE-Sign Loss	Embedding Initialization	mAP
Baseline	CLS1K			37.8
	CLS1K	✓		44.1
	CLS1K	✓	✓ BERT	53.5
	CLS1K	✓	✓ CLIP	54.7
DEFR/32	CLIP	✓	✓ CLIP	60.5
DEFR/16	CLIP	✓	✓ CLIP	65.6

proposed approach is effective on ResNet architectures as well.

4.3. Ablations

Several ablations were performed that focus on key components of DEFR: (a) the image backbone, (b) language embedding initialization and (c) the loss function.

From baseline to our solution: Table 4 shows the path from baseline to our DEFR. The baseline uses ImageNet pre-trained ViT-B/32 as backbone, random initialization for the classifier, and binary cross entropy loss. It achieves 37.8 mAP, much lower than the detection-assisted approaches (see Fig. 1). The LSE-Sign loss gains 6.3 mAP, and language embedding initialization adds another 9.4 mAP with BERT

Table 5. **Ablations on backbone pre-training and classifier initialization.** Numbers are mAP on the HICO dataset. *Random*: the default random initialization; *Embedding*: language embedding initialization using BERT or CLIP’s text encoder. The backbone is ViT-B/32 and fine-tuned with LSE-Sign loss.

Pre-training	Classifier Initialization		
	Random Initialization	BERT Embedding	CLIP Embedding
CLS1K	44.1	53.5 _(+9.4)	54.7 _(+10.6)
CLS21K	44.2	53.9 _(+9.7)	55.1 _(+10.9)
MLM	43.6	47.0 _(+3.4)	47.1 _(+3.5)
CLIP	36.8	51.0 _(+14.2)	60.5_(+23.7)

Table 6. **Binary Cross Entropy (BCE) Loss vs. LSE-Sign Loss** evaluated on HICO for four differently pre-trained models. The classifier is initialized with CLIP text embeddings.

Pre-training	Loss Function	
	BCE Loss	LSE-Sign Loss
	BERT Embedding Initialization	
CLS1K	50.6	53.5 _(+2.9)
CLS21K	51.0	53.9_(+2.9)
MLM	45.8	47.0 _(+1.2)
CLIP	44.4	51.0 _(+6.6)
	CLIP Embedding Initialization	
CLS1K	51.5	54.7 _(+3.2)
CLS21K	50.0	55.1 _(+5.1)
MLM	46.6	47.1 _(+0.5)
CLIP	57.9	60.5_(+2.6)

as the language model. Performance is further improved when the language model is jointly trained with the image backbone (to 60.5 mAP). It is worth mentioning that using the CLIP image backbone alone scores 34.2 mAP after fine-tuned on this task, much lower than our final solution. This demonstrates that our proposed approaches are effective and complementary mechanisms for HOI recognition.

Pre-training: We evaluate three different pre-training tasks for the backbone: (a) image classification (CLS1K/CLS21K), (b) masked language modeling (MLM), and (c) image-text contrastive learning (CLIP). Tab. 5 (the first column) shows the results for these pre-training tasks using random initialization in the linear classifier.

Classifier initialization: Table 5 compares the conventional random initialization and language embedding initialization. Clearly, the language embedding initialization with both language models provides a consistent improvement over all pre-training methods. It is worth noting that ImageNet pre-trained models (CLS1K/CLS21K) gain 10+ points from language embedding initialization.

Table 7. **Comparison between LSE-Sign loss and other loss functions** evaluated on HICO. We change the loss function of our best model with ViT-B/32 backbone and language embedding initialization. Weighted BCE is the binary cross entropy loss weighted by positive-negative-ratio per-class. Focal loss uses $\gamma=2$ and $\alpha=0.25$ as recommended in [31].

Loss function	mAP
Weighted BCE	54.7
BCE	57.9
Focal Loss	53.2
LSE-Sign Loss (ours)	60.5

Table 8. **Ablation of scalar γ** in Equation 1 evaluated on HICO dataset. The highest accuracy is achieved at $\gamma = 100$. The backbone is CLIP pre-trained ViT-B/32, the classifier is embedding initialized and LSE-Sign loss is used for fine-tuning.

γ	50	100	150	300	500
mAP	60.4	60.5	59.1	57.2	53.0

Loss function: The choice of loss function is vital in fine-tuning. Previous works use binary cross entropy (BCE) loss and treat HOI recognition as a set of binary classification problems. Table 6 shows that the proposed LSE-Sign loss outperforms BCE loss on all four differently pre-trained backbones. Table 7 compares LSE-Sign loss with other alternatives: binary cross entropy (BCE) loss, weighted BCE loss and focal loss [31] on our model. Weighted cross entropy loss is intended to impose a weight on the loss of each class so that each category is balanced among the positive samples and negative samples. It is not effective if the dataset is severely long-tailed. Focal loss [31] reduces the weight of the massive negative samples, but requires manual tuning of two hyper-parameters. Our LSE-Sign loss outperforms all others by a clear margin.

Scalar γ in Equation 1: LSE-Sign loss has a scalar γ which controls the magnitude of output per class s_i . Table 8 shows the accuracy under different values of γ . The maximum performance is achieved when $\gamma = 100$.

Backbone architecture: We train DEFR with different backbone architectures on HICO as in Tab. 9. We see all backbones outperform the current SOTA [28] by using our proposed methods.

4.4. HOI Detection

We find a properly trained HOI classification model can aid HOI detection effectively. We connect our frozen classification model with an off-the-shelf object detector to recognize the regional HOI. This detector-decoupled pipeline achieves SOTA performance in HOI detection, *without* train-

Table 9. **Ablation of backbone architecture** on HICO dataset. We apply our method on different ResNet and ViT backbone architectures. The pre-training of the backbone, language embedding initialization and loss function stay the same as our best model. DEFR method with ResNet-50 achieves 49.7 mAP, still outperforms the SOTA [28] by 2.6 mAP.

DEFR Backbone	mAP
ResNet-50	49.7
ResNet-101	53.6
ViT-B/32	60.5
ViT-B/16	65.6

ing DEFR on HICO-DET.

Dataset: We evaluate HOI detection performance on the HICO-DET dataset [4]. HICO-DET contains 117,871 annotated human-object pairs in the training set and 33,405 in the test set. HICO-DET and HICO share the same images and classes, but HICO-DET provides bounding box localized HOI annotations.

Evaluation metric: Following existing studies, we use the standard evaluation protocol [4]. Each positive prediction should have the correct HOI class together with a pair of human-object bounding boxes with IoU greater than 0.5 in reference to ground truth. Similar to the classification task, the average precision (AP) of each HOI class is separately computed and then averaged (mAP).

HOI detection pipeline: Leveraging the frozen DEFR classification model, we build an HOI detection pipeline without additional training, by simply connecting DEFR to an off-the-shelf object detector. The detected bounding boxes are converted to self-attention masks consumed by the last transformer layer in the ViT backbone (details in Appendix Sec. 7.1.1). Therefore, the CLS token in the last layer attends only to the region of interest specified by the pair of human and object bounding boxes. Object probabilities are multiplied to corresponding HOI classes.

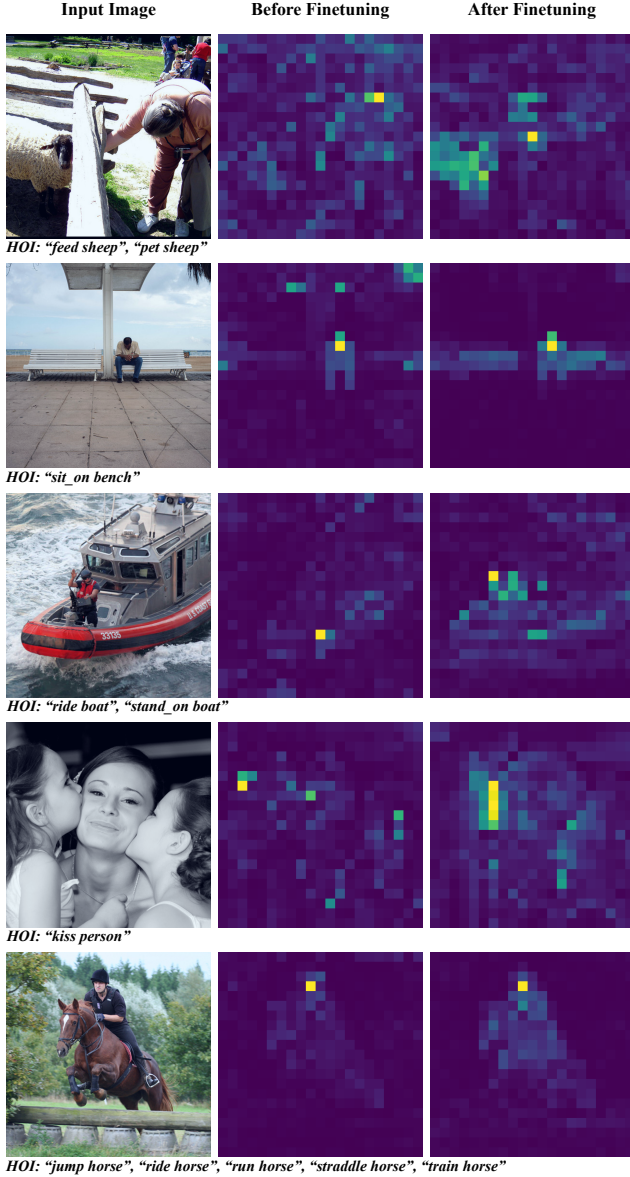
By doing so, we treat HOI detection as a special case of HOI classification, as in this task, DEFR classifies only the regional HOI. This pipeline requires no training on instance-level HOI annotation thanks to the properly learned feature in the classification task (see Fig. 5). Our approach has decoupled object detection from HOI recognition, hence more robust to the change of the detector (see Appendix Sec. 7.2).

4.5. Comparison on HOI detection

Table 10 compares the HOI detection performance on the HOI-DET dataset. Detectors used in this table are fine-tuned on the HICO-DET dataset. We use the detector from [10, 50].

Surprisingly, though DEFR is never trained under instance-level HOI supervision from HICO-DET, our

Figure 5. **Visualization of CLS attention map.** Left column: the input image. Middle column: the attention map of the CLIP pre-trained backbone without fine-tuning on HICO. Right column: the attention map of DEFR fine-tuned on HICO. After fine-tuning, the attention activates more on the HOI related objects and activates less elsewhere. The feature for HOI learned in the classification task aids HOI detection effectively.



method achieves 32.35 mAP, outperforming existing SOTA that are trained on HICO-DET. Specifically, on rare sets (classes with less than 10 samples), our method outperforms current studies by a clear margin (7.7 mAP), which is consistent with our image-level classification results. Different from previous works, our method showcases the strong correlation between HOI-Classification and HOI-Detection in a

Table 10. **Comparison of HOI detection performance** on HICO-DET [4]. The best performing method is in bold. Full: full set of 600 HOI classes, Rare: a subset of 138 HOI classes that have less than 10 training instances, Non-rare: the rest 462 classes. Methods [6, 21, 42] are end-to-end based on DETR [3]

Method	Backbone	Full	Rare	Non-rare
PPDM [30]	Hourglass-104	21.94	13.97	24.32
Bansal et al. [2]	ResNet-101	21.96	16.43	23.63
HOI-Trans [52]	ResNet-50	23.46	16.91	25.41
GG-Net [51]	Hourglass-104	23.47	16.48	25.60
VCL [18]	ResNet-50	23.63	17.21	25.55
ATL [19]	ResNet-50	23.81	17.43	25.72
DRG [10]	ResNet-50-FPN	24.53	19.47	26.04
HOTR [21]	ResNet-50	25.10	17.34	27.42
IDN [27]	ResNet-50	26.29	22.61	27.39
AS-Net [6]	ResNet-50	28.87	24.25	30.25
QPIC [42]	ResNet-50	29.07	21.85	31.23
SCG [50]	ResNet-50-FPN	31.33	24.72	33.31
CDN-B [49]	ResNet-50	31.78	27.55	33.05
CDN-L [49]	ResNet-101	32.07	27.19	33.53
Ours	ViT-B/16	32.35	33.45	32.02

very simple format that HOI-Detection is a special case of HOI-Classification with additional input of bounding boxes of human and objects. Note that these detected boxes is obtained from an object detector trained separately. This also demonstrates that our DEFR model is effective to learn a good representation for human-object interaction.

5. Discussion of Limitations

Although this paper reveals that language embedding initialization and LSE-Sign loss function play important roles for detection-free HOI recognition, we have not shown if these two components can work together with detection supervision to further push the state of the arts. In HOI detection, our performance on non-rare classes is lower than the state of the art. We will investigate these in future work.

6. Conclusion

In this paper, we decouple object and human keypoint detection from Human-Object Interaction (HOI) recognition. Our proposed detection-free method not only simplifies the pipeline, but also achieves higher accuracy than the detection-assisted counterparts. DEFR builds upon two findings. Firstly, we show that the structure of HOI classes can be effectively leveraged by using language embedding as classifier initialization. Secondly, we propose the LSE-Sign loss to facilitate multi-label learning on a long-tailed dataset. A combination of DEFR and an object detector achieves SOTA performance on HOI detection without additional fine-tuning. We hope that our work opens up a new direction for HOI recognition.

References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3686–3693, 2014. [5](#)
- [2] Ankan Bansal, Sai Saketh Rambhatla, Abhinav Shrivastava, and Rama Chellappa. Detecting human-object interactions via functional generalization. 2020. [8](#)
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In European Conference on Computer Vision, pages 213–229. Springer, 2020. [2](#), [8](#)
- [4] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In 2018 IEEE winter conference on applications of computer vision (wacv), pages 381–389. IEEE, 2018. [2](#), [7](#), [8](#), [11](#)
- [5] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiaxuan Wang, and Jia Deng. Hico: A benchmark for recognizing human-object interactions in images. In Proceedings of the IEEE International Conference on Computer Vision, pages 1017–1025, 2015. [1](#), [2](#), [3](#), [5](#)
- [6] Mingfei Chen, Yue Liao, Si Liu, Zhiyuan Chen, Fei Wang, and Chen Qian. Reformulating hoi detection as adaptive set prediction. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 9004–9013, June 2021. [2](#), [8](#)
- [7] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6163–6171, 2019. [2](#)
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020. [1](#), [3](#)
- [9] Hao-Shu Fang, Jinkun Cao, Yu-Wing Tai, and Cewu Lu. Pair-wise body-part attention for recognizing human-object interactions. In Proceedings of the European conference on computer vision (ECCV), pages 51–67, 2018. [1](#), [2](#), [3](#), [5](#), [6](#)
- [10] Chen Gao, Jiarui Xu, Yuliang Zou, and Jia-Bin Huang. Drg: Dual relation graph for human-object interaction detection. In European Conference on Computer Vision, pages 696–712. Springer, 2020. [2](#), [3](#), [7](#), [8](#), [11](#)
- [11] Chen Gao, Yuliang Zou, and Jia-Bin Huang. ican: Instance-centric attention network for human-object interaction detection. arXiv preprint arXiv:1808.10437, 2018. [2](#)
- [12] Rohit Girdhar and Deva Ramanan. Attentional pooling for action recognition. In Proceedings of the 31st International Conference on Neural Information Processing Systems, pages 33–44, 2017. [1](#), [2](#), [5](#), [6](#)
- [13] 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. [2](#)
- [14] Georgia Gkioxari, Ross Girshick, and Jitendra Malik. Contextual action recognition with r* cnn. In Proceedings of the IEEE international conference on computer vision, pages 1080–1088, 2015. [1](#), [2](#), [5](#), [6](#)
- [15] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Proceedings of the thirteenth international conference on artificial intelligence and statistics, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. [3](#)
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In Proceedings of the IEEE international conference on computer vision, pages 1026–1034, 2015. [3](#)
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016. [3](#)
- [18] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. Visual compositional learning for human-object interaction detection. 2020. [8](#)
- [19] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Affordance transfer learning for human-object interaction detection. In CVPR, 2021. [8](#)
- [20] Bumsoo Kim, Taeho Choi, Jaewoo Kang, and Hyunwoo J Kim. Uniondet: Union-level detector towards real-time human-object interaction detection. In European Conference on Computer Vision, pages 498–514. Springer, 2020. [2](#)
- [21] Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, and Hyunwoo J. Kim. Hotr: End-to-end human-object interaction detection with transformers. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 74–83, June 2021. [2](#), [8](#)
- [22] Dong-Jin Kim, Xiao Sun, Jinsoo Choi, Stephen Lin, and In So Kweon. Detecting human-object interactions with action co-occurrence priors. In European Conference on Computer Vision, pages 718–736. Springer, 2020. [2](#)
- [23] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. arXiv preprint arXiv:2102.03334, 2021. [5](#)
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. [5](#)
- [25] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. International journal of computer vision, 123(1):32–73, 2017. [2](#), [5](#)
- [26] Yikang Li, Wanli Ouyang, Bolei Zhou, Jianping Shi, Chao Zhang, and Xiaogang Wang. Factorizable net: an efficient subgraph-based framework for scene graph generation. In Proceedings of the European Conference on Computer Vision (ECCV), pages 335–351, 2018. [2](#)
- [27] Yong-Lu Li, Xinpeng Liu, Xiaoqian Wu, Yizhuo Li, and Cewu Lu. Hoi analysis: Integrating and decomposing human-object interaction. 2020. [8](#)

- [28] Yong-Lu Li, Liang Xu, Xinpeng Liu, Xijie Huang, Yue Xu, Mingyang Chen, Ze Ma, Shiyi Wang, Hao-Shu Fang, and Cewu Lu. Hake: Human activity knowledge engine. *arXiv preprint arXiv:1904.06539*, 2019. 1, 2, 3, 5, 6, 7
- [29] Yong-Lu Li, Liang Xu, Xinpeng Liu, Xijie Huang, Yue Xu, Shiyi Wang, Hao-Shu Fang, Ze Ma, Mingyang Chen, and Cewu Lu. Pastanet: Toward human activity knowledge engine. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 382–391, 2020. 1, 2, 3, 5, 6
- [30] 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. 2, 8
- [31] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 7
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 5
- [33] Yang Liu, Qingchao Chen, and Andrew Zisserman. Amplifying key cues for human-object-interaction detection. In *European Conference on Computer Vision*, pages 248–265. Springer, 2020. 2
- [34] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 5
- [35] Cewu Lu, Hao Su, Yonglu Li, Yongyi Lu, Li Yi, Chi-Keung Tang, and Leonidas J. Guibas. Beyond holistic object recognition: Enriching image understanding with part states. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1
- [36] Arun Mallya and Svetlana Lazebnik. Learning models for actions and person-object interactions with transfer to question answering. In *European Conference on Computer Vision*, pages 414–428. Springer, 2016. 1, 2, 5, 6, 11
- [37] Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. *Advances in neural information processing systems*, pages 570–576, 1998. 2
- [38] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24:1143–1151, 2011. 5
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 5
- [40] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018. 5
- [41] Yifan Sun, Changmao Cheng, Yuhan Zhang, Chi Zhang, Liang Zheng, Zhongdao Wang, and Yichen Wei. Circle loss: A unified perspective of pair similarity optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6398–6407, 2020. 3
- [42] Masato Tamura, Hiroki Ohashi, and Tomoaki Yoshinaga. Qpic: Query-based pairwise human-object interaction detection with image-wide contextual information. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10410–10419, June 2021. 2, 8
- [43] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Conference on Computer Vision and Pattern Recognition*, 2020. 2
- [44] Feng Wang, Xiang Xiang, Jian Cheng, and Alan Loddon Yuille. Normface: L2 hypersphere embedding for face verification. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1041–1049, 2017. 3
- [45] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5410–5419, 2017. 1, 2
- [46] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph r-cnn for scene graph generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–685, 2018. 2
- [47] Alireza Zareian, Haoxuan You, Zhecan Wang, and Shih-Fu Chang. Learning visual commonsense for robust scene graph generation. *arXiv preprint arXiv:2006.09623*, 2020. 2
- [48] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5831–5840, 2018. 2
- [49] Aixi Zhang, Yue Liao, Si Liu, Miao Lu, Yongliang Wang, Chen Gao, and Xiaobo Li. Mining the benefits of two-stage and one-stage hoi detection. *arXiv preprint arXiv:2108.05077*, 2021. 2, 8
- [50] Frederic Z. Zhang, Dylan Campbell, and Stephen Gould. Spatially conditioned graphs for detecting human-object interactions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13319–13327, October 2021. 2, 7, 8, 11
- [51] Xubin Zhong, Xian Qu, Changxing Ding, and Dacheng Tao. Glance and gaze: Inferring action-aware points for one-stage human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13234–13243, June 2021. 2, 8
- [52] Cheng Zou, Bohan Wang, Yue Hu, Junqi Liu, Qian Wu, Yu Zhao, Boxun Li, Chenguang Zhang, Chi Zhang, Yichen Wei, and Jian Sun. End-to-end human object interaction detection with hoi transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11825–11834, June 2021. 2, 8

7. Appendix

7.1. Ablations for HOI Detection Implementation

We compare two methods for using DEFR (ViT backbone) to classify the regional HOI in our HOI detection pipeline, given the bounding boxes for a human-object pair:

- **ave-pool**: use the *average feature* of the patches inside the bounding boxes for classification.
- **mask-n**: extract a *regional CLS* token for classification by applying self-attention masks in the last n transformer layers. The mask prevents CLS attending to patches outside the bounding boxes.

7.1.1 Implementation Details

The self-attention mask Φ is applied in the *Attention* [36] function of a transformer layer:

$$\text{Attention}(Q, K, V) = \text{softmax}(\Phi + \frac{QK^T}{\sqrt{d_k}})V \quad (4)$$

Φ is a binary mask converted from the bounding boxes. $\Phi_{i,j}$ equals $-\infty$ if i is the CLS token and j a patch outside the given bounding boxes, and 0 otherwise. d_k is the dimension of Q, K and V . The *Attention* function is used in each head as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (5)$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

where W^O, W_i^Q, W_i^K, W_i^V are trainable projection matrices as in [36].

For *ave-pool*, Eq. (4) is changed to:

$$\text{Attention}(Q, K, V) = \text{softmax}(\Phi)V \quad (6)$$

where $\Phi_{i,j}$ equals $-\infty$ if i is the CLS token and j either the CLS token or a patch outside the given bounding boxes, and 0 otherwise. Then, the CLS output becomes the mean of tokens inside the area of interest and is used for classification.

Both methods do not change the dimension of the backbone output, which makes it possible to reuse the pre-trained classification head directly.

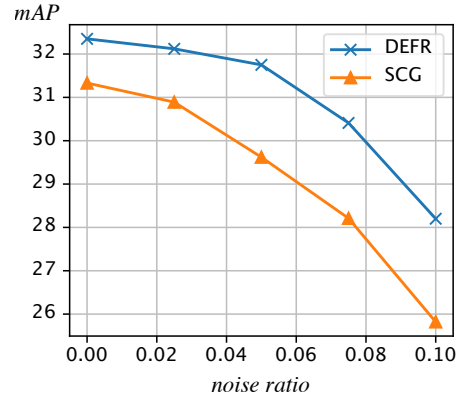
7.1.2 Results

Tab. 11 shows that applying self-attention mask in the final layer achieves the best performance, since it is most similar to training. Specifically, simply averaging the features inside the boxes (ave-pool) performs much worse, because the informative feature is diluted by the background inside the bounding boxes, besides, the pre-trained Q, K matrices are not utilized in that layer. Using the self-attention effectively is essential to leverage the full potential of DEFR for HOI detection.

Table 11. **HOI detection** evaluated on HICO-DET [4]. DEFR with the ViT-B/16 backbone is connected with the same detector used in [10, 50]. Regional feature extraction methods *ave-pool* and *mask-n* are compared, where for *mask-n*, $n=[1, \dots, 6]$ are used.

Method	Full	Rare	Non-rare
DEFR (ave-pool)	15.32	11.64	16.42
DEFR (mask-6)	31.27	32.45	30.91
DEFR (mask-5)	31.29	32.52	30.92
DEFR (mask-4)	31.37	32.48	31.04
DEFR (mask-3)	31.41	32.60	31.06
DEFR (mask-2)	31.59	32.53	31.31
DEFR (mask-1)	32.35	33.45	32.02

Figure 6. **Analysis of robustness to detector change** on the HICO-DET dataset. We add random noise to the bounding box coordinates at various levels (x-axis) and compare the performance of our method and the two-stage state-of-the-art SCG [50]. Our model shows less performance degradation.



7.2. HOI Detection Robustness Analysis

Our detector-decoupled design allows our method to work with any off-the-shelf object detector without additional training. Since DEFR is not jointly trained with the detector, our method can be more robust to the change of the detector. In this experiment, we add noise of various levels to the bounding box coordinates to simulate using different detectors, and compare the performance degradation with SCG [50].

Specifically, Gaussian noise is added to the bounding box coordinates as:

$$\begin{aligned} X' &= X + \mathcal{N}(0, r \cdot W) \\ Y' &= Y + \mathcal{N}(0, r \cdot H) \end{aligned} \quad (7)$$

where X, Y are original bounding box coordinates, W, H the widths and heights of the original bounding boxes, and X', Y' the noise-disturbed box coordinates used for HOI detection. The result in Fig. 6 shows that our detector-decoupled method has more robust performance.