# Improving Human-Object Interaction Detection via Virtual Image Learning

Shuman Fang
Key Laboratory of Multimedia
Trusted Perception and Efficient
Computing, Ministry of Education of
China, Xiamen University
fangshuman@stu.xmu.edu.cn

Shuai Liu
Key Laboratory of Multimedia
Trusted Perception and Efficient
Computing, Ministry of Education of
China, Xiamen University
luckyliu@stu.xmu.edu.cn

Jie Li
Key Laboratory of Multimedia
Trusted Perception and Efficient
Computing, Ministry of Education of
China, Xiamen University
lijie.32@outlook.com

Guannan Jiang
Intelligent Manufacturing
Department, Contemporary Amperex
Technology Co. Limited (CATL)
jianggn@catl.com

Xianming Lin*
Key Laboratory of Multimedia
Trusted Perception and Efficient
Computing, Ministry of Education of
China, Xiamen University
linxm@xmu.edu.cn

Rongrong Ji
Key Laboratory of Multimedia
Trusted Perception and Efficient
Computing, Ministry of Education of
China, Xiamen University
rrji@xmu.edu.cn

## ABSTRACT

Human-Object Interaction (HOI) detection aims to understand the interactions between humans and objects, which plays a curtail role in high-level semantic understanding tasks. However, most works pursue designing better architectures to learn overall features more efficiently, while ignoring the long-tail nature of interaction-object pair categories. In this paper, we propose to alleviate the impact of such an unbalanced distribution via *Virtual Image Leaning* (VIL). Firstly, a novel label-to-image approach, *Multiple Steps Image Creation* (MUSIC), is proposed to create a high-quality dataset that has a consistent distribution with real images. In this stage, virtual images are generated based on prompts with specific characterizations and selected by multi-filtering processes. Secondly, we use both virtual and real images to train the model with the teacher-student framework. Considering the initial labels of some virtual images are inaccurate and inadequate, we devise an *Adaptive Matching-and-Filtering* (AMF) module to construct pseudo-labels. Our method is independent of the internal structure of HOI detectors, so it can be combined with off-the-shelf methods by training merely 10 additional epochs. With the assistance of our method, multiple methods obtain significant improvements, and new state-of-the-art results are achieved on two benchmarks.

## CCS CONCEPTS

• **Computing methodologies → Object detection**; **Activity recognition and understanding**.

*Corresponding author.

## KEYWORDS

Human-Object Interaction Detection, Long-tail Distribution

## 1 INTRODUCTION

HOI detection is to comprehend the interactive relationships between humans and objects, which can be denoted as a set of triplets ⟨*human, object, interaction*⟩. It has attracted the interest of many researchers due to its strong correlation with other vision tasks. It not only contributes to other high-level semantic understanding tasks [2, 22, 34, 41] but also benefits from basic vision tasks [3, 14, 15, 42]. However, the long-tail distributions for interactions and objects are common in the dataset. The combinatorial nature of HOIs further exacerbates the number gap between rare and non-rare categories [11–13, 36, 40]. Models trained on such datasets only fit well in the common categories, while ignoring the rare ones.

To address the long-tail issue, re-sampling and re-weighting are designed to make detectors focus on tailed categories [11, 36, 44]. But they underperform due to insufficient diverse features. Some methods [13, 17, 40] involve latent linguistic embeddings of rare categories to augment feature space, which suffers from lacking visual representations yet. To get diverse visual features, ATL [12] expands training images by introducing extra off-the-shelf datasets. Nevertheless, gathering large-scale relative images is challenging. Generating virtual images becomes a straightforward idea. Considering Stable Diffision (SD) [28] has succeeded in generating high-quality images, we pursue a more suitable way for HOI detection to augment real datasets with it.

In this work, we propose a training framework termed *Virtual Image Leaning* (VIL). The paradigm comparison of existing methods and ours are shown in Figure 1. Firstly, to ensure the consistency
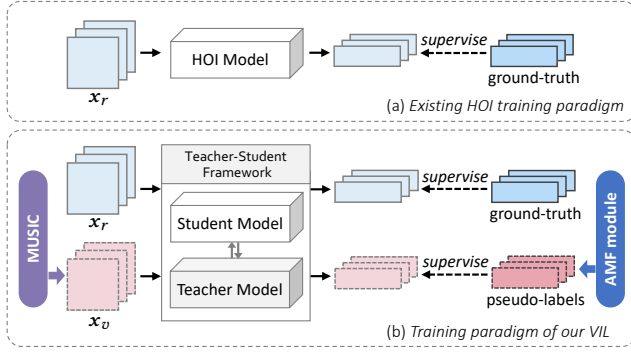
**Figure 1: Comparison for the HOI training paradigm. (a) Most existing HOI detectors receive real images $x_r$ to predict HOI triplets, supervised by ground truth of $x_r$. (b) Our VIL is trained with both virtual images $x_v$ and real images $x_r$, with the supervision of pseudo-labels and ground truth. During this, virtual images are generated from MUSIC, and pseudo-labels are constructed by the AMF module.**

of virtual and real images, we propose a novel label-to-image generation approach, named ***Mul**tiple **S**teps **I**mage **C**reation* (MUSIC). Specially, we first create natural linguistic descriptions based on object and interaction categories. Then, add characterizations for humans to improve the quality of virtual images and constraints for scenes to ensure the authenticity of interactions. The retouched textual descriptions are received by Stable Diffusion [28] to acquire virtual images. To eliminate the uncertainty and instability in the generation process, we discard low-confidence virtual images with several filter mechanisms, including scene similarity, instance existence, and human-object interactiveness. With this, an unlimited number of virtual images with annotations can be obtained.

Considering the initial annotations of virtual images are noisy and incomplete, we have to correct inaccurate bounding boxes and complement the missing HOI triplets. Hence, we refer to the teacher-student framework and propose *Adaptive Matching-and-Filtering* (AMF) module to create pseudo-labels for the virtual images. To improve the accuracy of bounding boxes, the AMF module computes matching costs adaptively to avoid interference from noisy boxes. And adaptive thresholds are used to pick up high-confidence predictions to recall HOI triplets in the images. The student model learns from both virtual and real images, supervised by the pseudo-labels and the ground-truth labels. Meanwhile, the teacher model receives the student model's feedback to improve the quality of pseudo-labels in the next iteration.

Our proposed method is simple, general, and orthogonal to existing methods. By combining with our method, almost all off-the-shelf HOI detectors can be improved with only additional 10 training epochs. To evaluate its efficiency, we conduct extensive experiments with multiple off-the-shelf HOI detectors on two widely used datasets. All experimented methods achieve relative gains and new state-of-the-art results are obtained on two datasets. Ablation studies verify the contribution of each part of VIL. Visualization experiments of virtual images and pseudo-labels explain the source of efficiency of our proposed MUSIC and AMF module. Concretely, we summarize our contribution as follows:

- We propose VIL, a model-agnostic framework, to boost the detection performance of existing methods.
- We devise a label-to-image generation approach named MUSIC, which can generate high-quality virtual images that have consistent distribution with real images.
- We design AMF, a pseudo-label generation module, to correct and supplement the initial labels of virtual images.
- By combining multiple methods with ours, the performance gains of all methods and the new state-of-the-art results indicate the efficacy of VIL.

## 2 RELATED WORK

### Category Bias Solution in HOI Detection

The performance of many existing HOI methods [6, 19, 20, 35, 50] is limited by the long-tail issue. Methods proposed to address it can be categorized into three streams: re-sampling, re-weighting, and data space extension.

VCL [11] emphasizes the long-tail issue for the first time. Given an input image, VCL randomly samples another one and permutes interaction-object pair in these two images to obtain new combinations, which somehow alleviates the long-tail dilemma. ODM [36] proposes another sampling strategy by alternately performing write-in and read-out stages. The write-in stage dynamically updates the memory with rare categories' features, whereas the read-out part seeks to sample the far-distance features.

A dynamic re-weighting mechanism is proposed by CDN [44], which amplifies the weight of rare categories during extra training epochs with a relatively small learning rate. However, both re-sampling and re-weighting suffer from overfitting existing features.

To relieve overfitting to the existing rare representations, an effective solution is to expand the available space of data or features. FCL [13], based on word embeddings of interaction-object pair categories, generates virtual features from Gaussian noise to enrich feature space. Instead of using latent features which lack visual representations, ATL [12] directly extends the original dataset with additional datasets [21, 30]. Learning from such rich and varied knowledge, the model can get rid of unbalanced distribution. However, such additional datasets are hard to obtain and still suffer from the limited number of images. Consequently, what we pursue is to generate large-scale virtual images with distributions that are consistent with the original ones.

### Data Augmentation Based on Stable Diffusion

The development of the diffusion model [10] has enabled current study, Stable Diffusion (SD) [28], to produce high-quality images with remarkable progress. It allow general conditioning inputs (*e.g.*, text) to synthesize images. Considering its success in image generation, numerous researchers have investigated its application in data augmentation [1, 7, 9, 29, 33, 51]. Almost existing SD-based augmentation works focus on the instance-level. They commit to improving portraying realistic objects. However, images generated by such approaches are unsuitable for HOI detection. In addition to photo-realistic humans and objects, the synthetic images utilized for HOI identiìcation should also represent speciìc scene information and plausible interaction behaviors. To this end, we propose MUSIC to address the inadequacies based on them.
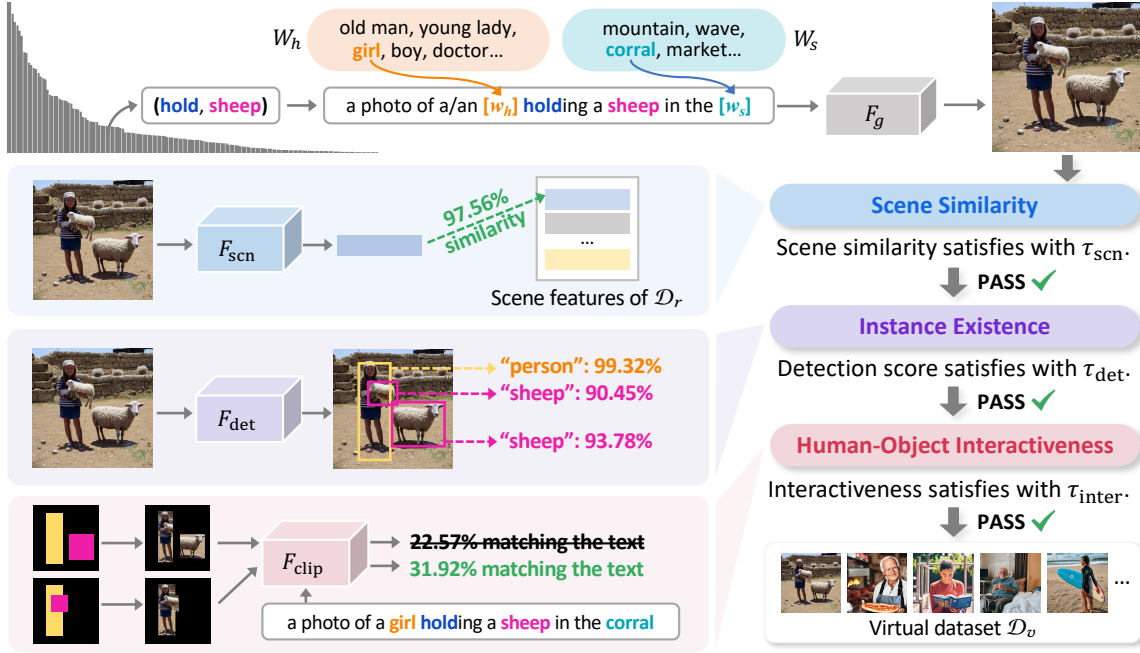
**Figure 2: Overview of MUSIC. Given a category $(c_a, c_o)$, MUSIC firstly expands it into a plain linguistic sentence and refines it by adding human and scene depictions. A text-to-image model $F_g$ receives the refined sentence as the prompt to generate a virtual image. Then the image will be evaluated by Scene Similarity, Instance Existence, and Human-Object Interactiveness to determine whether to keep it or not.**

## 3   VIRTUAL IMAGE LEARNING

### 3.1   Overview

Our proposed VIL first generates a virtual dataset $\mathcal{D}_v$ based on the long-tail distribution, where a label-to-image generation approach, MUSIC (described in Section 3.2), is introduced. Then, the teacher-student framework is adopted to train on both virtual dataset $\mathcal{D}_v$ and real dataset $\mathcal{D}_r$. During this stage, we design an AMF module (presented in Section 3.3) to generate pseudo-labels for virtual images. The student model learns knowledge supervised with the pseudo-labels of virtual images and the ground truth of real images. Meanwhile, the teacher model is updated by the student model to provide more reliable pseudo-labels in the next iteration.

### 3.2   Virtual Image Generation

Due to the excellent image generation ability of Stable Diffusion [28], recent researchers have an appetite for expanding datasets with it [7, 51]. This text-to-image model can output the images corresponding to the given text prompt. In this work, we also adopt it to generate virtual images based on HOI categories. In particular, given an interaction category $c_a$ and an object category $c_o$, we first expand the pair $(c_a, c_o)$ into a description sentence: "a photo of a person $[c_a$-ing] a/an $[c_o]$". Then, Stable Diffusion, denoted as $F_g$, receives the sentence as the prompt to generate the corresponding image. We name this approach ***Direct Image Creation*** (DIC). However, the virtual images created by DIC can not achieve satisfactory quality. As Figure 4 shown, we analyze the failure cases come from the following reasons: 1) for humans, only parts of human-body

are visible; 2) for scenes, the background can not provide indicative information; 3) for interactiveness, there may not be interactions between humans and objects. Based on these, the MUSIC is proposed to handle the issues. It consists of text refinement, Scene Similarity, Instance Existence, and Human-Object Interactiveness. The overall procedure of MUSIC is depicted in Figure 2.

*3.2.1 Text Refinement.* As mentioned above, the virtual images generated by plain descriptions only contain human body parts, lacking of human pose features. We argue that human body parts alone can not reflect human appearance characteristics. To close to the input prompt, model $F_g$ should generate an image with an integrated person by portraying his face or clothing. Consequently, we prepare a word set $W_h$. The words in it are representative of human characteristics, like age, gender, occupation, *etc.* MUSIC refines the label-based sentence by replacing "person" with a specific human characteristic word $w_h$ randomly sampled from $W_h$. To make the scene indicative, MUSIC also adds scene description into the prompt sentence. In particular, another word set $W_s$ is also needed, which is composed of scene categories in Places365 [49]. For each $(c_a, c_o)$, we count all possible scene categories, represented as $W_s^{(c_a, c_o)} \subseteq W_s$. And the scene description $w_s$ is sampled from $W_s^{(c_a, c_o)}$. To sum up, the plain sentence is polished by adding human and scene depiction, which can be unified as "a photo of a/an $[w_h]$ $[c_a$-ing] a/an $[c_o]$ in the $[w_s]$". We denote the refined sentence as $t$ and send it into $F_g$ to generate the virtual image $x_v$. Next, multiple filtering stages will evaluate for $x_v$ to determine whether to keep it or drop it.

Shuman Fang, Shuai Liu, Jie Li, Guannan Jiang, Xianming Lin, and Rongrong Ji

*3.2.2 Scene Similarity.* We argue that the scene in virtual image $x_v$ should be similar to that in real images, so we compute the cosine similarity of scene features for evaluation. In this stage, a scene prediction model $F_{scn}$ [49] is introduced to extra scene features. The scene quality score $s_{scn}$ of $x_v$ can be expressed as the following formula:

$$s_{scn} = \max_{x_r \in \mathcal{D}_r} \frac{\left(F_{scn}(x_v)\right)^\top \left(F_{scn}(x_r)\right)}{\|F_{scn}(x_v)\|_2 \cdot \|F_{scn}(x_r)\|_2}, \tag{1}$$

where $\mathcal{D}_r$ is the real dataset, and $\|\cdot\|_2$ is $\ell_2$ norm. The virtual image $x_v$ can be determined as passable to the scene similarity filtering only if $s_{scn}$ is greater than the threshold $\tau_{scn}$. Otherwise, it will be abandoned.

*3.2.3 Instance Existence.* To avoid target humans or objects being too over-occluded to be detected, $x_v$ needs to be verified by Instance Existence filtering. Specifically, $x_v$ is sent to a pre-trained object detector $F_{det}$ [3] to obtain a set of predictions. MUSIC selects all predictions predicted as "person" or $c_o$ with confidence greater than $\tau_{det}$. Then split them into two candidate sets $B_h$ and $B_o$. If the lengths of the two candidate sets are both greater than 0, it is considered that $x_v$ clearly contains a human and the target object. If not, such an image will be discarded.

*3.2.4 Human-Object Interactiveness.* In this turn, we need to verify whether $x_v$ accurately expresses the semantic information of the refined sentence $t$. In particular, MUSIC enumerates the bounding boxes in $B_h$ and $B_o$ to combine them into a set of human-object pairs. Each pair can be denoted as $(b_h, b_o)$. The bounding box can be further represented as $b_\xi = \left[x_1^\xi, y_1^\xi, x_2^\xi, y_2^\xi\right]^\top$, where $\xi \in \{h, o\}$. We mask the pixels in human or object regions to obtain:

$$x_{mask} = \mathbf{M} \odot x_v,$$
$$\mathbf{M}_{(i,j)} = \begin{cases} 1, \text{if } i \in [x_1^\xi, x_2^\xi] \wedge j \in [y_1^\xi, y_2^\xi] \\ 0, \text{otherwise} \end{cases}. \tag{2}$$

By traversing all pairs, we can get a masked image set $X_{mask}$.

After that, MUSIC uses a visual-linguistic model $F_{clip}$ [26] to compute the semantic similarity between each masked image and the refined sentence $t$. The maximum similarity will be regarded as the score $s_{inter}$ of this filtering step:

$$s_{inter} = \max_{x_{mask} \in X_{mask}} F_{clip}(x_{mask}, t). \tag{3}$$

The virtual image $x_v$ can pass through this step if $s_{inter}$ satisfies the threshold $\tau_{inter}$. And the human-object pair $(\hat{b}_h, \hat{b}_o)$ that corresponds to the maximum similarity is served as the annotation bounding boxes for $x_v$. Otherwise, $x_v$ will be dropped.

The virtual dataset $\mathcal{D}_v$ consists of all virtual images passing through all three filtering processes, formulated by:

$$\mathcal{D}_v = \left\{\left(x_v^i, y_v^i\right)\right\}_{i=1}^{N_v},$$
$$y_v = \left(c_a, c_o, \hat{b}_o, \hat{b}_h\right) \in \{1, \cdots, C_a\} \times \{1, \cdots, C_o\} \times \mathbb{R}^4 \times \mathbb{R}^4, \tag{4}$$

where $N_v$ is the size of the virtual dataset, $C_a$ and $C_o$ are the numbers of interaction categories and object categories in the real dataset $\mathcal{D}_r$. With the help of MUSIC, the distribution of the virtual dataset is consistent with that of the real dataset. In Figure 4, we show some examples generated by MUSIC.

## 3.3 Virtual Image Training

We argue that the initial annotations of virtual images still show deficiencies. On the one hand, some of them contain incorrect bounding boxes, which means there do not exist interactions between the located human and object. On the other hand, the number of HOI triplets in initial annotations is insufficient. Annotation in each image only includes one HOI triplet since we adopt a one-label-to-one-image strategy to generate virtual. But in fact, there are multiple triplets for each image. Hence, a pseudo-labels module should be designed to correct and complement the initial annotations so that pseudo-labels can supervise the learning of the virtual images better. To end this, we propose the AMF module to generate high-quality pseudo-labels and introduce the teacher-student framework to train the model by learning from both virtual images and real images.

*3.3.1 Pseudo-Labels Generation.* As mentioned above, our proposed AMF module aims to 1) correct the wrong bounding boxes and 2) supplement the missing HOI triplets.

To achieve the former objective, we adaptively compute the matching cost to find the most similar prediction. In particular, given virtual image $x_v$, its corresponding initial annotation is $y_v = \left(c_a, c_o, \hat{b}_o, \hat{b}_h\right)$. By sending $x_v$ into the teacher model, the predictions consist of the following four parts: the human bounding boxes $\{\tilde{b}_h^i \mid \tilde{b}_h^i \in \mathbb{R}^4\}_{i=1}^N$, the object bounding boxes $\{\tilde{b}_o^i \mid \tilde{b}_o^i \in \mathbb{R}^4\}_{i=1}^N$, the probability of object classes $\{\tilde{s}_o^i \mid \tilde{s}_o^i \in [0,1]^{C_o+1}\}_{i=1}^N$, and the probability of interaction classes $\{\tilde{s}_a^i \mid \tilde{s}_a^i \in [0,1]^{C_a}\}_{i=1}^N$, where $N$ is the size of the prediction set. For the $i$-th prediction, we formulate the matching cost of classification $H_{cls}^i$ as the following:

$$H_{cls}^i = H_a^i + H_o^i,$$
$$H_a^i = -\frac{1}{2}\left(\tilde{s}_a^i[c_a] + \frac{1}{N-1}\sum_{k \in \{1 \cdots N\} \setminus \{c_a\}} 1 - \tilde{s}_a^i[k]\right), \tag{5}$$
$$H_o^i = -\tilde{s}_o^i[c_o],$$

where $[\cdot]$ means index operation. And the localization cost $H_{loc}^i$ can be formulated like the following:

$$H_{loc}^i = H_{reg}^i + H_{iou}^i,$$
$$H_{reg}^i = \max\left\{\left\|\tilde{b}_h^i - \hat{b}_h\right\|_1, \left\|\tilde{b}_o^i - \hat{b}_o\right\|_1\right\}, \tag{6}$$
$$H_{iou}^i = \max\left\{1 - GIoU(\tilde{b}_h^i, \hat{b}_h), 1 - GIoU(\tilde{b}_o^i, \hat{b}_o)\right\},$$

where $GIoU(\cdot)$ is the generalized IoU [27]. Considering the bounding boxes may be inaccurate while the classes are absolutely correct, we compute the overall matching cost by adaptively dropping the localization part, that is:

$$H^i = \begin{cases} H_{cls}^i, \text{if } \min_k H_{cls}^k + H_{loc}^k > 0 \\ H_{cls}^i + H_{loc}^i, \text{otherwise} \end{cases}. \tag{7}$$

With this, we can adaptively find the nearest prediction by searching for $\omega = \arg\min H^i$ with the Hungarian algorithm [18]. To correct the initial bounding boxes, we replace $\hat{b}_h$ with $\tilde{b}_h^\omega$ and $\hat{b}_o$ with $\tilde{b}_o^\omega$, respectively. Also, we set the $\tilde{s}_a^\omega[c_a]$ as infinite to guarantee this prediction can be picked up in the high-confidence filtering.

For the latter objective, *i.e.*, the number of HOI triplets in the initial annotation is insufficient, we select high-confidence predictions to supplement it. Considering the confidence gap among different HOI detectors, we suggest using an adapt threshold to select predictions. Firstly, we estimate the average number of human-object pairs in each of virtual images and denote the number as $\kappa$. When the teacher model is introduced, we calculate its prediction scores for all virtual images. And we select $(\kappa \times N_v)$-th highest score as the threshold $\tau_{\text{bin}}$. All predictions with interaction confidence higher than $\tau_{\text{bin}}$ are picked up, where pair-wise NMS [44] is introduced to remove duplicate predictions. By defining a score binarization function $f_{\text{bin}}(s) = \lceil s - \tau_{\text{bin}} \rceil$, we can get the final pseudo-labels:

$$\tilde{y}_v = \left\{ \left( f_{\text{bin}}(\tilde{s}_a^i), c_o, \tilde{b}_h^i, \tilde{b}_o^i \right) \mid \max \tilde{s}_a^i > \tau_{\text{bin}} \right\}. \tag{8}$$

With high-quality pseudo-labels, the student model can learn from virtual images much more effectively.

*3.3.2 Teacher-Student Framework.* Inspired by Omni-DETR [38], we apply strong augmentation $T^s$ and weak augmentation $T^w$ to both virtual images $x_v$ and real images $x_r$. And for virtual images $x_v$, we additionally design a random padding augmentation to avoid the model overfitting the large-area bounding boxes, which can be formulated as:

$$T^p(x_v, y_v) = \begin{cases} \text{RandomPad}(x_v, y_v), \text{if } S_{\hat{b}_h} > \frac{1}{2} S_{x_v}, p \le 0.5 \\ (x_v, y_v), \text{otherwise} \end{cases}, \tag{9}$$

where $\text{RandomPad}(\cdot)$ represents the random padding function, the variable $p \sim U(0, 1)$ is introduced to control the transformation rate, and $S_{\hat{b}_h}$ and $S_{x_v}$ are areas of the human bounding box and the virtual image. Thus, for each virtual data $(x_v, y_v) \in \mathcal{D}_v$, we can get $(x_v^{ps}, y_v^{ps})$ and $(x_v^{pw}, y_v^{pw})$, where the superscript $ps$ means apply $T^p$ and $T^s$ successively, and similarly for $pw$. For each real image $(x_r, y_r) \in \mathcal{D}_r$, they are transformed into $(x_r^s, y_r^s)$ and $(x_r^w, y_r^w)$.

In the training stage, $x_v^{pw}$ will be sent into the teacher model $\mathcal{F}_t(x, \theta_t)$ to predict pseudo-labels $\tilde{y}_v^{pw}$, which is described in Section 3.3.1. Then we restore $\tilde{y}_v^{pw}$ by the inverse transformation of $T^w$ and apply the strong augmentation the same as $x_v^{ps}$ to get the transformed pseudo-labels $\tilde{y}_v^{ps}$. Finally, the student model $\mathcal{F}_s(x, \theta_s)$ is trained with $x_v^{ps}$, $x_r^s$, and $x_r^w$ under the supervision of $\tilde{y}_v^{ps}$, $y_r^s$, and $y_r^w$, respectively. The total loss function to optimize it can be represented as:

$$\begin{aligned} L_{\text{total}} = \sum_{x_v \in \mathcal{D}_v} L_{\text{hoi}}\left( \mathcal{F}_s(x_v^{ps}), \tilde{y}_v^{ps} \right) \\ + \sum_{x_r \in \mathcal{D}_r} L_{\text{hoi}}\left( \mathcal{F}_s(x_r^s), y_r^s \right) + L_{\text{hoi}}\left( \mathcal{F}_s(x_r^w), y_r^w \right), \end{aligned} \tag{10}$$

where $L_{\text{hoi}}(\cdot, \cdot)$ denotes the loss function utilized in the off-the-shelf HOI detectors.

Meanwhile, the teacher model should be updated by the exponential moving average (EMA) [32] from the student model:

$$\theta_t \leftarrow \alpha \theta_t + (1 - \alpha) \theta_s, \tag{11}$$

where $\alpha$ is a hyperparameter empirically set to a number close to 1 to keep the robustness in the teacher model.

**Table 1: Comparison with existing long-tail methods by combining QPIC [31] on the HICO-Det dataset. The content in parentheses indicates the performance improvements.**

| Method | Full | Rare | Non-rare |
|---|---|---|---|
| QPIC | 29.07 | 21.85 | 31.23 |
| QPIC + ATL [12] | 28.87 (−0.20) | 21.67 (−0.23) | 31.03 (−0.20) |
| QPIC + ODM [36] | 29.26 (+0.19) | 22.07 (+0.22) | 31.41 (+0.18) |
| QPIC + CDN [44] | 29.40 (+0.33) | 21.96 (+0.11) | 31.63 (+0.40) |
| QPIC + VIL (ours) | **30.54 (+1.47)** | **23.34 (+1.49)** | **32.69 (+1.46)** |
| GEN-VLKT | 33.75 | 29.25 | 35.10 |
| GEN-VLKT + ODM [36] | 33.82 (+0.07) | 29.59 (+0.34) | 35.08 (−0.02) |
| GEN-VLKT + CDN [44] | 33.87 (+0.12) | 29.35 (+0.10) | 35.22 (+0.12) |
| GEN-VLKT + VIL (ours) | **34.21 (+0.46)** | **30.58 (+1.33)** | **35.30 (+0.20)** |

## 4 EXPERIMENT

### 4.1 Experiment Setup

*4.1.1 Dataset and Metric.* We follow previous works to evaluate performance on two public benchmarks, *i.e.*, HICO-Det [4] and V-COCO [8]. We compute the mean average precision (mAP) to report experimental results. A prediction is a true positive if the predicted human and object boxes have IoUs larger than 0.5 with the corresponding ground truth and the predicted HOI category is also correct.

In the HICO-Det dataset, there are 38,118 images for training and 9,658 for testing. Images are annotated with 80 object classes and 117 action classes. The HOI category is defined as an interaction-object pair (*e.g.*, "eat pizza"). There are two settings for evaluation: Default and Known Object. The Default setting requires evaluating all images, while Known Object only tests images containing the target object class. In each setting, HICO-Det also provides three evaluation sets, *i.e.*, Full, Rare, and Non-Rare, which are divided based on the frequency of categories.

V-COCO dataset contains 5,400 images in the trainval set and 4,946 in the test set. Images are annotated with 80 object classes and 29 interaction classes. Four interaction classes (*i.e.*, stand, walk, run, and smile) are neglected for evaluation since they are not associated with semantic roles. The HOI category is defined as an interaction class. We report role mAP in two scenarios, where scenario 1 needs to predict the cases in which humans interact with no objects while scenario 2 ignores these cases.

*4.1.2 Implementation Details.* In the virtual image generation process, we set the thresholds $\tau_{\text{scn}}$, $\tau_{\text{det}}$, and $\tau_{\text{inter}}$ as 0.9, 0.9, and 0.3, respectively. For the HICO-Det dataset, the number of virtual images for each interaction-object pair category is set to 40 per rare category and 10 per non-rare category. For the V-COCO dataset, we first split all object-action pair categories into two groups: the minority with less than 10 instances in the train set and the majority with 10 or more. And we generate 30 and 15 virtual images for each minority and majority category, respectively.

During the training stage, $\alpha$ of EMA is set to 0.9996, the same as [32, 38]. Considering that our VIL is orthogonal with most existing methods, we conduct experiments by applying ours to them. Hence, we set the rest hyperparameters the same as the methods to be combined. We initialize the student and teacher model by loading parameters from the corresponding pre-trained method

Shuman Fang, Shuai Liu, Jie Li, Guannan Jiang, Xianming Lin, and Rongrong Ji

**Table 2: Performance improvement on both HICO-Det and V-COCO datasets. Each letter in the Feature column stands for A: Appearance/Visual feature, S: Spatial features, L: Linguistic feature of label semantic embeddings, and P: Human pose feature. \* signifies results reproduced with the official implementation codes. The performance improvements in the Full and Rare sets are marked with RED and BLUE, respectively.**

| Method | Backbone | Feature | HICO-Det | | | | | | V-COCO | |
| | | | Default | | | Known Object | | | | |
| | | | Full | Rare | Non-rare | Full | Rare | Non-rare | Scenario 1 | Scenario 2 |
| DRG [5] | ResNet-50-FPN | A+S+P+L | 19.26 | 17.74 | 19.71 | 23.40 | 21.75 | 23.89 | 51.0 | - |
| VSGNet [35] | ResNet-152 | A+S | 19.80 | 16.05 | 20.91 | - | - | - | 51.8 | - |
| PPDM [19] | Hourglass-104 | A | 21.73 | 13.78 | 24.10 | 24.58 | 16.65 | 26.84 | - | - |
| GG-Net [48] | Hourglass-104 | A | 23.47 | 16.48 | 25.60 | 27.36 | 20.23 | 29.48 | - | - |
| SCG [45] | ResNet-50-FPN | A+S | 31.33 | 24.72 | 33.31 | 34.37 | 27.18 | 36.52 | 54.2 | 60.9 |
| PhraseHOI [23] | ResNet-50 | A+L | 29.29 | 22.03 | 31.46 | 31.97 | 23.99 | 34.36 | 57.4 | - |
| CPC [24] | ResNet-50 | A | 29.63 | 23.14 | 31.57 | - | - | - | 63.1 | 65.4 |
| SSRT [16] | ResNet-50 | A+L | 30.36 | 25.42 | 31.83 | - | - | - | 63.7 | 65.9 |
| HQM [47] | ResNet-50 | A | 31.34 | 26.54 | 32.78 | - | - | - | 63.6 | - |
| UPT [46] | ResNet-50 | A+S | 31.66 | 25.94 | 33.36 | 35.05 | 29.27 | 36.77 | 59.0 | 64.5 |
| CDN [44] | ResNet-50 | A | 31.78 | 27.55 | 33.05 | 34.53 | 29.73 | 35.96 | 61.7 | 63.8 |
| SDT [37] | ResNet-50 | A | 32.45 | 28.09 | 33.75 | 35.95 | 31.30 | 37.34 | 60.3 | 65.7 |
| QPIC [31] | ResNet-50 | A | 29.07 | 21.85 | 31.23 | 31.68 | 24.14 | 33.93 | 58.8 | 61.0 |
| +*VIL (ours)* | ResNet-50 | A | 30.54 (+1.47) | 23.34 | 32.69 | 33.24 (+1.56) | 25.13 | 35.66 | 59.4 (+0.6) | 61.9 (+0.9) |
| OCN [43] | ResNet-50 | A+L | 30.91 | 25.56 | 32.51 | 33.68* | 28.27* | 35.30* | 64.2 | 66.3 |
| +*VIL (ours)* | ResNet-50 | A+L | 31.99 (+1.08) | 26.67 | 33.58 | 34.75 (+1.07) | 29.49 | 36.32 | 64.9 (+0.7) | 67.0 (+0.7) |
| DOQ [25] | ResNet-50 | A+L | 31.55 | 26.75 | 32.99 | 34.11* | 29.25* | 35.55* | 63.5 | 65.9* |
| +*VIL (ours)* | ResNet-50 | A+L | 32.40 (+0.85) | 27.95 | 33.73 | 34.99 (+0.88) | 30.19 | 36.42 | 64.3 (+0.8) | 67.1 (+1.2) |
| DisTR [50] | ResNet-50 | A | 31.93* | 27.26* | 33.32* | 34.62* | 29.53* | 36.14* | 66.4* | 68.6* |
| +*VIL (ours)* | ResNet-50 | A | 32.84 (+0.91) | 28.04 | 34.27 | 35.63 (+1.01) | 30.53 | 37.15 | **67.6** (+1.2) | **69.9** (+1.3) |
| GEN-VLKT [20] | ResNet-50 | A+L | 33.75 | 29.25 | 35.10 | 36.78 | 32.75 | 37.99 | 64.6* | 66.8* |
| +*VIL (ours)* | ResNet-50 | A+L | **34.21** (+0.46) | **30.58** | **35.30** | **37.67** (+0.89) | **34.88** | **38.50** | 65.3 (+0.7) | 67.7 (+0.9) |

and train the framework with 10 epochs. The learning rate of the backbone and the other parts are set to the same as those after decay in the to-be-combined methods. Considering virtual images are introduced due to category bias, we freeze the parameters except for classification heads when training with virtual images to prevent overfitting.

## 4.2 Comparison with Long-tail Methods

In Table 1, we compare our VIL with existing long-tail methods with QPIC [31] as the baseline. Note that "+CDN" means only applying the re-weighting technique in CDN. From the table, only ATL brings the performance decline to QPIC. Since ATL is proposed based on traditional two-stage HOI methods, it needs the multi-stream structure to fuse features from additional datasets flexibly. Such a design makes it incompatible with transformer-based methods. Instead, our method is model-agnostic and suitable for almost all HOI methods. Compared with the other methods, the improvement from our VIL also outstands those from others by a large margin on all sets. We owe this to the various visual features provided by virtual images, which is the shortage of re-sampling [36] and re-weighting [44] techniques. Moreover, by taking the SOTA method, GEN-VLKT [20], as baseline, our method still outperforms existing works. We also visualize the improvements from CDN and our VIL in Figure 3. From the figure, the more rare the category is, the more improvement the baseline obtains. And our improvement under

rare categories is far more than CDN. We conclude that our method can address long-tail problem much more effectively.

## 4.3 Improvement on Existing Works

To evaluate the efficiency and generalization of our VIL, we select five representative methods to conduct combination experiments, which are one classic method QPIC [31], two relatively recent works that introduced extra knowledge OCN [43] and DOQ [25], and two SOTA methods GEN-VLKT [20] and DisTR [50].

We first conduct experiments on the HICO-Det test set and report results in Table 2. All methods achieve performance gains after combining with VIL. For example, QPIC with the help of VIL, improved by 1.47 mAP and 1.56 mAP under the Default and the Known Object setting, surpassing many recent methods[16, 23, 24, 36]. OCN and DOQ get 1.08 and 0.85 mAP gains, respectively, outperforming almost all existing works. Similarly, the improvement of DisTR achieves 0.91 mAP. The state-of-the-art method on HICO-Det, GEN-VLKT, can also be further enhanced to achieve a new SOTA result. Note that in the Rare set under the two different settings, it acquires considerable margins of 1.33 and 2.13 mAP, with relative improvements achieving 4.55% and 6.50%. We attribute the improvements on the Rare set to our long-tail-based design, which enables the model to pay more attention to features in rare categories.

For performance on the V-COCO test set, we still combine our VIL with the five works, which are QPIC, DOQ, OCN, GEN-VLKT,
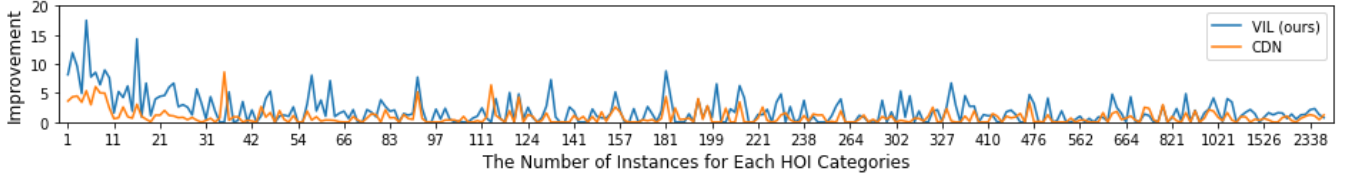
**Figure 3: Comparison improvements with our VIL and CDN [44] by taking QPIC [31] as baseline method on the HICO-Det dataset under the Default setting. We sort the categories by their frequency. The horizontal axis is the number of instances in the train set, and the vertical axis is the mean improvements of categories with the same number of instances.**

**Table 3: Ablation experiments for MUSIC approach (depicted in Section 3.2) on the V-COCO test set.**

| #Row | Text | Scene | Instance | Interactiveness | Scenario 1 | Scenario 2 |
|------|------|-------|----------|-----------------|------------|------------|
| **Ours** | ✓ | ✓ | ✓ | ✓ | **59.4** | **61.9** |
| 1 | ✓ | ✓ | ✓ | | 59.2 (−0.2) | 61.6 (−0.3) |
| 2 | ✓ | ✓ | | | 59.8 (−0.5) | 61.3 (−0.6) |
| 3 | ✓ | | | | 58.7 (−0.7) | 61.1 (−0.8) |
| 4 | | | | | 58.2 (−1.2) | 60.7 (−1.2) |

**Table 4: Ablation experiments for AMF module (depicted in Section 3.3) on V-COCO test set.**

| #Row | Matching | Filtering | Scenario 1 | Scenario 2 |
|------|----------|-----------|------------|------------|
| **Ours** | ✓ | ✓ | **59.4** | **61.9** |
| 1 | ✓ | | 59.2 (−0.2) | 61.6 (−0.3) |
| 2 | | | 58.8 (−0.6) | 61.4 (−0.4) |

and DisTR. The results are also reported in Table 2, where all methods are relatively improved. Specifically, the early work, QPIC, attains 0.6 and 0.9 mAP in Scenario 1 and 2, making it competitive with those latter works [23, 46]. The other three recent works, DOQ, OCN, and GEN-VLKT, are enhanced by 0.8, 0.7, and 0.7 mAP, respectively. DisTR, the SOTA method on the V-COCO dataset, earns more significant improvement by 1.2 and 1.3 mAP in Scenario 1 and 2. It achieves outstanding results of 67.3 and 69.7 mAP in the two scenarios, refreshing the SOTA results.

## 4.4 Ablation Study

Firstly, to verify the efficiency of each part of MUSIC, we conduct ablation experiments for it and report the results in Table 3, where the header "Text", "Scene", "Instance", and "Interactiveness" are corresponding the part depicted in Section 3.2.1, 3.2.2, 3.2.3, and 3.2.4. As shown results, the performance degradations verify the necessity of each part. Also, the performance plummets dramatically when getting rid of Text Refinement (row 4). We argue that the characterizations of humans play a quite crucial role in virtual images, which is in line with intuition. And the three last parts serve to validate the elements that need to be portrayed.

Then, we prove the effectiveness of the AMF module by removing "Matching" and "Filtering" one by one. The results are shown in Table 4. Without Filtering (row 1), learning one virtual image is supervised by only one HOI annotation. The insufficient annotations lead to a 0.2 mAP decline. By removing the Matching part (row 2), the further decline illustrates the detrimental impact of noisy

bounding box annotations. Moreover, the performance degradation in the two rows also illustrates that the AMF module can construct more reliable pseudo-labels to guide the student model.

## 4.5 Qualitative Results

*4.5.1 Visualization for Virtual Images.* The ablation study has demonstrated that the virtual images generated by DIC may impair the performance of existing methods, whereas our MUSIC does the opposite. To demonstrate this more intuitively, we visualize the virtual images generated by these two approaches for comparison in Figure 4, where images at the first row are results from DIC and the rest are from MUSIC. It is clear that the DIC images primarily lack human bodies, meaningful sceneries, or interactive human-object pairs. For example, all shown images only contain human hands. Among those, the hand in column 1 is severely occluded, and that in column 2 is illegible. And images in columns 2, 5, and 8 miss the background information, while such information has been proved to be crucial to the understanding of interaction [25, 31, 39, 45]. Additionally, for the image in column 6, the interaction "lay" between the human and the couch is not distinguished enough. In contrast, the images generated via MUSIC can successfully address the aforementioned issues. All images can perfectly convey the semantic information of the specified categories, which benefits from the multiple filtering stages to guarantee the quality of virtual images. Moreover, the images have diversity in the appearance characteristics of the humans, the posture of the human bodies, and the scene where the interaction occurs. In particular, MUSIC can depict two different cases, "cut hair with scissors" and "cut paper with scissors", based on the category "(cut_instr, scissors)". For the category "(lay, couch)", various lying postures are constructed. Apart from these, for the categories "(read, book)", "(talk, cell_phone)", and "(hit_instr, tennis_racket)", MUSIC can provide a variety of reasonable scenes. We owe these to the introduced textual descriptions that direct model $F_g$ to create various virtual images.

*4.5.2 Visualization for Pseudo-Labels.* Considering the noisy initial annotations created during the generation stage, AMF module is proposed to construct more reliable pseudo-labels. Many initial annotations are with erroneous bounding boxes or lack sufficient HOI triplets. We visualize some examples in Figure 5. In the first row, the initial annotation localizes the "person" incorrectly. It is the boy that looking at a pizza, not the man in the background. In this case, we expect the pseudo-label to correct the human bounding box. At the beginning epochs, the pseudo-labels successfully fix the bounding box of the human but conflate a bunch of pizzas. As

Shuman Fang, Shuai Liu, Jie Li, Guannan Jiang, Xianming Lin, and Rongrong Ji



**Figure 4: Comparison of virtual image samples generated by DIC (first row) and MUSIC (last two rows). The interaction-object pair categories used to generate virtual images are marked at the top of images, where the interaction and object classes are in BLUE and PINK, respectively.**



**Figure 5: Visualization for the change process of the pseudo-labels. We pick up some cases where the initial annotation created during the generation stage is inaccurate and show them on the leftmost with RED. The pseudo-labels during training epochs are on the right, where the interactive pairs belonging to the same group are drawn with the same color.**

the training progresses, with the human bounding box keeping precise, the bounding box of the pizza instance gradually moves closer to the correct direction, and finally locates exactly. In the second row, the initial annotation ignores the interaction between "person" and "baseball". As the pseudo-labels show, the interactive relationship is dug out at the 1st epoch. But this relationship is lost again due to the instability of predictions. As the teacher model is continuously updated, this initially missed interactive triplet can be pointed out stably after the 6th epoch. We emphasize the necessity of pseudo-labels, which can provide better and richer supervised information for the student model.

## 5 CONCLUSION

In this paper, we track the problem of long-tail the dilemma in HOI detection. We propose a novel and general framework termed *Virtual Image Leaning* (VIL) to enhance existing HOI detectors. In particular, to generate a large-scale high-quality virtual dataset, we design a *Multiple Steps Image Creation* (MUSIC) approach. Given an interaction-object pair category, MUSIC expands it to a plain sentence, polishes it with specific descriptions, and evaluates the generated image by multiple filtering steps. In the training stage, *Adaptive Matching-and-Filtering* (AMF) module is adopted to denoise and supplement the initial annotations of virtual images. And the obtained pseudo-labels, along with the groud-truth, supervise the model learning knowledge from virtual and real images. By combining five representative methods, we evaluate the effectiveness and generalization of our VIL on two public datasets. The results demonstrate all methods make significant progress with our method, and new state-of-the-art performances emerge.

# REFERENCES

[1] Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J Fleet. 2023. Synthetic data from diffusion models improves imagenet classification. *arXiv preprint arXiv:2304.08466* (2023).

[2] Matteo Bruni, Tiberio Uricchio, Lorenzo Seidenari, and Alberto Del Bimbo. 2016. Do Textual Descriptions Help Action Recognition?. In *ACM International Conference on Multimedia*.

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European Conference on Computer Vision*.

[4] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. 2018. Learning to detect human-object interactions. In *Workshop on Applications of Computer Vision*.

[5] Chen Gao, Jiarui Xu, Yuliang Zou, and Jia-Bin Huang. 2020. Drg: Dual relation graph for human-object interaction detection. In *European Conference on Computer Vision*.

[6] Chen Gao, Yuliang Zou, and Jia-Bin Huang. 2018. iCAN: Instance-centric attention network for human-object interaction detection. In *British Machine Vision Conference*.

[7] Yunhao Ge, Jiashu Xu, Brian Nlong Zhao, Laurent Itti, and Vibhav Vineet. 2022. Dall-e for detection: Language-driven context image synthesis for object detection. *arXiv preprint arXiv:2206.09592* (2022).

[8] Saurabh Gupta and Jitendra Malik. 2015. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474* (2015).

[9] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. 2023. Is synthetic data from generative models ready for image recognition?. In *International Conference on Computer Vision*.

[10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*.

[11] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. 2020. Visual compositional learning for human-object interaction detection. In *European Conference on Computer Vision*.

[12] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. 2021. Affordance transfer learning for human-object interaction detection. In *Computer Vision and Pattern Recognition*.

[13] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. 2021. Detecting human-object interaction via fabricated compositional learning. In *Computer Vision and Pattern Recognition*.

[14] Jie Hu, Liujuan Cao, Yao Lu, ShengChuan Zhang, Yan Wang, Ke Li, Feiyue Huang, Ling Shao, and Rongrong Ji. 2021. Istr: End-to-end instance segmentation with transformers. *arXiv preprint arXiv:2105.00637* (2021).

[15] Jie Hu, Linyan Huang, Tianhe Ren, Shengchuan Zhang, Rongrong Ji, and Liujuan Cao. 2023. You Only Segment Once: Towards Real-Time Panoptic Segmentation. In *Computer Vision and Pattern Recognition*.

[16] ASM Iftekhar, Hao Chen, Kaustav Kundu, Xinyu Li, Joseph Tighe, and Davide Modolo. 2022. What to look at and where: Semantic and spatial refined transformer for detecting human-object interactions. In *Computer Vision and Pattern Recognition*.

[17] Zhong Ji, Xiyao Liu, Yanwei Pang, and Xuelong Li. 2020. SGAP-Net: Semantic-guided attentive prototypes network for few-shot human-object interaction recognition. In *Association for the Advancement of Artificial Intelligence*.

[18] Harold W Kuhn. 1955. The Hungarian method for the assignment problem. *Naval research logistics quarterly* (1955).

[19] Yue Liao, Si Liu, Fei Wang, Yanjie Chen, Chen Qian, and Jiashi Feng. 2020. Ppdm: Parallel point detection and matching for real-time human-object interaction detection. In *Computer Vision and Pattern Recognition*.

[20] Yue Liao, Aixi Zhang, Miao Lu, Yongliang Wang, Xiaobo Li, and Si Liu. 2022. Gen-vlkt: Simplify association and enhance interaction understanding for hoi detection. In *Computer Vision and Pattern Recognition*.

[21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*.

[22] Fei Liu, Jing Liu, Zhiwei Fang, Richang Hong, and Hanqing Lu. 2020. Visual question answering with dense inter-and intra-modality interactions. *IEEE Transactions on Multimedia* (2020).

[23] Xinpeng Liu, Yong-Lu Li, Xiaoqian Wu, Yu-Wing Tai, Cewu Lu, and Chi-Keung Tang. 2022. Interactiveness Field in Human-Object Interactions. In *Computer Vision and Pattern Recognition*.

[24] Jihwan Park, SeungJun Lee, Hwan Heo, Hyeong Kyu Choi, and Hyunwoo J Kim. 2022. Consistency learning via decoding path augmentation for transformers in human object interaction detection. In *Computer Vision and Pattern Recognition*.

[25] Xian Qu, Changxing Ding, Xingao Li, Xubin Zhong, and Dacheng Tao. 2022. Distillation using oracle queries for transformer-based human-object interaction detection. In *Computer Vision and Pattern Recognition*.

[26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.

[27] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. 2019. Generalized intersection over union: A metric and a loss for bounding box regression. In *Computer Vision and Pattern Recognition*.

[28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Computer Vision and Pattern Recognition*.

[29] Mert Bulent Sariyildiz, Karteek Alahari, Diane Larlus, and Yannis Kalantidis. 2023. Fake it till you make it: Learning transferable representations from synthetic ImageNet clones. In *Computer Vision and Pattern Recognition*.

[30] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. 2019. Objects365: A large-scale, high-quality dataset for object detection. In *International Conference on Computer Vision*.

[31] Masato Tamura, Hiroki Ohashi, and Tomoaki Yoshinaga. 2021. QPIC: Query-Based Pairwise Human-Object Interaction Detection with Image-Wide Contextual Information. In *Computer Vision and Pattern Recognition*.

[32] Antti Tarvainen and Harri Valpola. 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in Neural Information Processing Systems*.

[33] Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. 2023. Effective data augmentation with diffusion models. *arXiv preprint arXiv:2302.07944* (2023).

[34] TJ Tsai, Andreas Stolcke, and Malcolm Slaney. 2015. A study of multimodal addressee detection in human-human-computer interaction. *IEEE Transactions on Multimedia* (2015).

[35] Oytun Ulutan, ASM Iftekhar, and Bangalore S Manjunath. 2020. Vsgnet: Spatial attention network for detecting human object interactions using graph convolutions. In *Computer Vision and Pattern Recognition*.

[36] Guangzhi Wang, Yangyang Guo, Yongkang Wong, and Mohan Kankanhalli. 2022. Chairs can be stood on: Overcoming object bias in human-object interaction detection. In *European Conference on Computer Vision*.

[37] Guangzhi Wang, Yangyang Guo, Yongkang Wong, and Mohan Kankanhalli. 2022. Distance Matters in human-object interaction detection. In *ACM International Conference on Multimedia*.

[38] Pei Wang, Zhaowei Cai, Hao Yang, Gurumurthy Swaminathan, Nuno Vasconcelos, Bernt Schiele, and Stefano Soatto. 2022. Omni-DETR: Omni-supervised object detection with transformers. In *Computer Vision and Pattern Recognition*.

[39] Tiancai Wang, Rao Muhammad Anwer, Muhammad Haris Khan, Fahad Shahbaz Khan, Yanwei Pang, Ling Shao, and Jorma Laaksonen. 2019. Deep contextual attention for human-object interaction detection. In *International Conference on Computer Vision*.

[40] Bingjie Xu, Yongkang Wong, Junnan Li, Qi Zhao, and Mohan S Kankanhalli. 2019. Learning to detect human-object interactions with knowledge. In *Computer Vision and Pattern Recognition*.

[41] Rui Yan, Peng Huang, Xiangbo Shu, Junhao Zhang, Yonghua Pan, and Jinhui Tang. 2022. Look less think more: Rethinking compositional action recognition. In *Proceedings of the 30th ACM International Conference on Multimedia*.

[42] Wei Yang, Shuang Li, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. 2017. Learning feature pyramids for human pose estimation. In *International Conference on Computer Vision*.

[43] Hangjie Yuan, Mang Wang, Dong Ni, and Liangpeng Xu. 2022. Detecting human-object interactions with object-guided cross-modal calibrated semantics. In *Association for the Advancement of Artificial Intelligence*.

[44] Aixi Zhang, Yue Liao, Si Liu, Miao Lu, Yongliang Wang, Chen Gao, and Xiaobo Li. 2021. Mining the benefits of two-stage and one-stage hoi detection. *Advances in Neural Information Processing Systems*.

[45] Frederic Z Zhang, Dylan Campbell, and Stephen Gould. 2021. Spatially conditioned graphs for detecting human-object interactions. In *International Conference on Computer Vision*.

[46] Frederic Z Zhang, Dylan Campbell, and Stephen Gould. 2022. Efficient two-stage detection of human-object interactions with a novel unary-pairwise transformer. In *Computer Vision and Pattern Recognition*.

[47] Xubin Zhong, Changxing Ding, Zijian Li, and Shaoli Huang. 2022. Towards hard-positive query mining for detr-based human-object interaction detection. In *European Conference on Computer Vision*.

[48] Xubin Zhong, Xian Qu, Changxing Ding, and Dacheng Tao. 2021. Glance and Gaze: Inferring Action-aware Points for One-Stage Human-Object Interaction Detection. In *Computer Vision and Pattern Recognition*.

[49] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017).

[50] Desen Zhou, Zhichao Liu, Jian Wang, Leshan Wang, Tao Hu, Errui Ding, and Jingdong Wang. 2022. Human-object interaction detection via disentangled transformer. In *Computer Vision and Pattern Recognition*.

[51] Shangchen Zhou, Kelvin Chan, Chongyi Li, and Chen Change Loy. 2022. Towards robust blind face restoration with codebook lookup transformer. In *Advances in Neural Information Processing Systems*.