

# Weakly Supervised Human-Object Interaction Detection in Video via Contrastive Spatiotemporal Regions

Shuang Li<sup>1\*</sup> Yilun Du<sup>1</sup> Antonio Torralba<sup>1</sup> Josef Sivic<sup>2</sup> Bryan Russell<sup>3</sup>  
<sup>1</sup>MIT <sup>2</sup>CIIRC CTU <sup>3</sup>Adobe

<https://shuangli-project.github.io/weakly-supervised-human-object-detection-video>

## Abstract

We introduce the task of weakly supervised learning for detecting human and object interactions in videos. Our task poses unique challenges as a system does not know what types of human-object interactions are present in a video or the actual spatiotemporal location of the human and the object. To address these challenges, we introduce a contrastive weakly supervised training loss that aims to jointly associate spatiotemporal regions in a video with an action and object vocabulary and encourage temporal continuity of the visual appearance of moving objects as a form of self-supervision. To train our model, we introduce a dataset comprising over 6.5k videos with human-object interaction annotations that have been semi-automatically curated from sentence captions associated with the videos. We demonstrate improved performance over weakly supervised baselines adapted to our task on our video dataset.

## 1. Introduction

In this paper, we study the problem of weakly supervised human-object interaction detection in videos. Given a video sequence, as illustrated in Figure 1, a system must correctly identify and localize the person and interacted object (“bike”) in the scene, in addition to identifying the action (“washing”) taken by the human, for the duration of the interaction in the video without bounding box supervision. While there has been an impressive progress in learning visual-language representations [38, 22, 28] from hundreds of millions of captioned images or videos recently, the learnt representations focus on classifying or retrieving entire images or videos given a language query. Our task is more challenging as it requires the models to correctly detect both the human and object bounding boxes in multiple frames of the video.

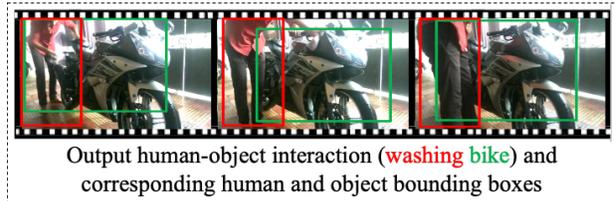


Figure 1: We seek to detect human-object interactions in videos. In this example, our system is able to detect “human washing bike” in the given video. Our approach learns to detect such interactions in a weakly supervised fashion, *i.e.*, without requiring bounding box annotations at training time. (Video credit: Dude Chennai [6])

Human-object interaction detection has been primarily studied in the context of still images [3, 4, 15, 20, 53, 54, 37, 56, 66, 52, 46]. However, they are naturally temporal events that take place over a period of time. Interactions such as “drinking” or “pushing” occur between a human and an object over time, making videos a natural modality for studying this problem.

Existing video-based methods primarily rely on strong bounding box supervision and having access to a fully annotated video dataset. However, relying on strong supervision has significant drawbacks. First, exhaustively annotating the spatial location of objects in a video is time consuming given the large number of frames in a video. Second, scaling to the large number of possible interactions and obtaining a sufficient number of ground truth bounding boxes is challenging due to the potentially open vocabulary of objects and actions and the combinatorial nature of human-object interactions. Third, interactions typically follow a long-tailed distribution, with common human-object interactions occurring much more frequently than others [47, 20]. While supervised learning usually prefers common interactions, a robust human-object interaction detection system should instead perform equally well on both common and rare interactions.

In this work, we seek to leverage videos with verb and noun phrase annotations derived from natural language sentence captions to learn to detect human-object interactions in videos in a weakly supervised manner. Such an approach

<sup>2</sup>Czech Institute of Informatics, Robotics and Cybernetics at the Czech Technical University in Prague

\*Work done at Adobe Research during SL’s summer internship

is advantageous as obtaining video-level annotations is significantly less costly than bounding boxes in videos. Leveraging such data makes it possible to scale training to a larger number of videos and vocabulary of objects and actions.

Our task is challenging as we do not know the correspondence between the verb-object queries and spatiotemporal regions in the training videos. A system must learn to establish these correspondences without spatial bounding box supervision. We thus propose a contrastive loss over spatiotemporal regions for detecting human-object interactions in videos. Our loss jointly associates candidate spatiotemporal regions with an action and object vocabulary in a weakly supervised manner and leverages cues about the temporal continuity of objects in motion as a form of self-supervision. Such a formulation allows us to deal with an open vocabulary of language queries which is especially desirable in human-object interaction, due to the high prevalence of rare and unseen action and object combinations.

Our paper has three main contributions: (1) We present an approach that integrates spatiotemporal information for humans and objects for weakly supervised human-object interaction detection in videos. Our approach does not require manual bounding box annotations. (2) We present a contrastive loss over spatiotemporal regions that leverages weak verb-object supervision from video captions and self-supervision from temporal continuity in video. It allows detecting rare and unseen human-object interactions in a zero-shot manner. (3) We introduce a new dataset of over 6.5k videos to evaluate human-object interaction in videos. We demonstrate improved performance over weakly supervised baselines adapted to our task. The dataset is made public to facilitate further research<sup>1</sup>.

## 2. Related work

Closest to our approach is work in modelling video and natural language, visual relationship detection, and human-object interaction detection.

**Video and natural language.** Prior work has looked at jointly modeling video and natural language for tasks, such as captioning [25], movie question answering [48], and short clip retrieval [40, 57]. More relevant are works that aim to more finely “ground” or align natural language in videos. Examples include retrieving moments from untrimmed videos [12, 18], learning from video with aligned instructions [29, 28], and alignment of natural language with (spatio-) temporal regions in a video [21]. Natural language poses hard challenges due to large open vocabulary and complex interactions due to composition.

**Visual relationship detection.** Previous work, *e.g.*, [2, 8, 14, 15, 16, 26, 35, 36, 41, 56, 59, 61], has studied detecting subject-predicate-object visual relations in single

still images. This line of work has been extended to video with strong supervision [43, 50]. Closest to our approach is work on weakly supervised visual relationship detection [34, 62, 55, 58] where a model is trained to use triplet annotation available at the image level. Different than us, Peyre *et al.* [34] leverage a fixed vocabulary of pre-trained object detectors and learn relations with a discriminative clustering model. Peyre *et al.* [33] model open language but in the strongly supervised setting and for still images.

**Human-object interaction detection.** Human-object interaction detection [4, 53, 54, 37, 66, 52] is a kind of human-centric relation detection. HOI is an essential research topic for deeper scene understanding. Several datasets, such as HICO-DET [3] and V-COCO [15], have been proposed for this domain. [45, 56, 20] formulate the novel HOI detection as a zero-shot learning problem. However, these methods are based on still images and have difficulties in detecting dynamic human-object interactions. They either rely on the bounding box annotations or pre-trained object detectors which has been show perform badly in videos [11].

## 3. Learning contrastive spatiotemporal regions

We address the problem of detecting human-object interactions (HOIs) in videos in a weakly supervised manner. As obtaining ground truth bounding boxes for supervised learning is expensive and time consuming, we seek to learn from a collection of videos where only verb-object phrase annotations are provided for the entire video clip during training. We thus propose a weakly supervised framework that incorporates both spatial and temporal information to detect HOIs in videos. The overall training setup is illustrated in Figure 2. Given a video clip and a verb-object query, for each frame, we first extract a bank of features. The features include those for the verb-object query, frame, and human/object regions in the clip. This bank of features passes through a region attention module that outputs two features for the frame – an attended human feature and an attended object feature that focus attention on regions that are more relevant to the verb-object query. These features, along with the verb-object feature and object region features from the other frames are passed into our weakly supervised contrastive loss.

### 3.1. Weakly supervised contrastive loss

Learning from language labels in a weakly supervised manner is challenging as a system must automatically identify and associate video spatiotemporal regions with the provided phrase annotations. Moreover, HOIs typically follow a long-tailed distribution. Applying the often used classification loss will not suffice as it requires a fixed vocabulary with similar number of samples for each class. Furthermore, a classification loss maximizes the probabil-

<sup>1</sup>Code and dataset are available at <https://shuangli-project.github.io/VHICO-Dataset>.

ity of the correct class while suppressing all other classes, which does not allow for less common or unseen objects and verbs. Finally, words with similar meanings are not explicitly mapped to nearby locations in the feature space with a classification loss.

To address these issues, we introduce a contrastive spatiotemporal loss for learning a shared visual-language embedding, as shown in Figure 3. Our loss leverages the phrase annotations associated with each training video and cues about the temporal continuity of objects in motion. Our training loss incorporates three insights. First, we learn to map the visual representation for the likely human and object regions to the corresponding embedded representation of the input verb-object queries and contrast against embedded representations of other non-relevant words in the vocabulary. Second, we encourage spatiotemporal regions to be temporally consistent in the video. Third, we apply the contrastive loss in our model, enabling it to detect new unseen human-object interactions during testing.

We build on the contrastive loss [5, 17, 19], which aims to encourage positive pairs of unit-length features to be close (measured by dot product) and negative pairs to be far in the feature space,

$$\mathcal{L}_C(f, f', \{f_n\}_{n=1}^N) = -f^T f' + \log \sum_{n=1}^N \exp(f^T f_n), \quad (1)$$

where  $f$  is an anchor feature,  $f'$  is a positive feature and  $\{f_n\}_{n=1}^N$  are  $N$  negative features. We propose a weakly supervised language-embedding alignment loss to align the spatiotemporal regions to the input verb-object query and a self-supervised temporal contrastive loss to encourage temporal continuity of the object regions based on Equation (1).

**Weakly supervised language-embedding alignment loss.** Given a video frame  $I_t$ , we extract its human and object region proposal features,  $f_t^h$  and  $f_t^o$ . Let  $e$  be a language-embedding feature for the ground truth verb-object label of the input video. We seek to align relevant human/object regions to the ground truth verb-object label. Since only the frame-level (or video-level) verb-object label is available, we also seek to learn a global human/object feature in each frame that contrasts against a negative set of language-embedding features  $\mathcal{E}$  covering the vocabulary not including the ground truth verb-object label.

To perform the alignment, we propose a region attention module that computes an attention score  $\sigma_{t,i}^h$  and  $\sigma_{t,i}^o$  for each human and object region proposal, respectively, to measure their relevance to the verb-object query. We obtain an attended human feature  $\Phi_t^h$  by aggregating the human region features  $f_t^h$  in frame  $I_t$  as a weighted average over their attention scores  $\sigma_t^h$ ,

$$\Phi_t^h = \sum_{i=1}^{N_h} \sigma_{t,i}^h f_{t,i}^h, \quad (2)$$

where  $N_h$  is the number of candidate human regions. The attended object feature  $\Phi_t^o$  has a similar form. The feature

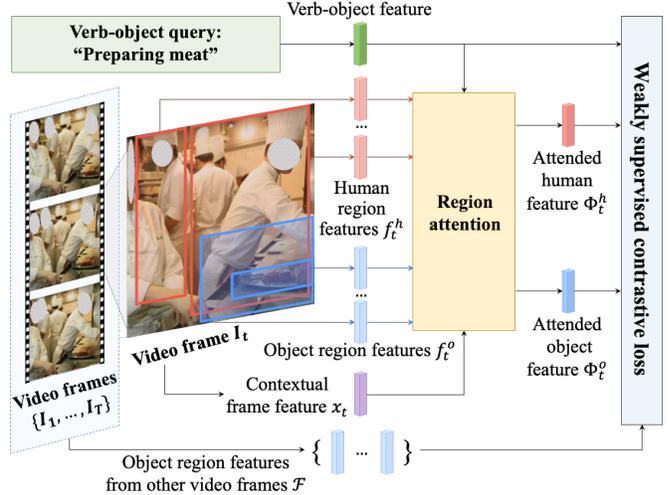


Figure 2: **Training overview.** Given a video clip and a verb-object query, for each frame, we first extract its human and object region features. The human/object features are aggregated in a region attention module to attend to regions that are more relevant to the query. The attended human feature, attended object feature, the feature of verb-object query, and object region features from other frames are used to compute our weakly supervised contrastive loss. (Video credit: The Best Gallery Craft [7])

attention “softly” selects a small number of candidate human/object regions as targets, with higher-scoring regions contributing more to the attended feature.

We define the language-embedding alignment loss  $\mathcal{L}_L$  as the alignment of the attended features in a frame to the target label while contrasting against the verb or object negative feature set. Following the general expression of the contrastive loss in Equation (1), we define the language-embedding alignment loss in frame  $I_t$  as a summation of contrastive losses given attended human/object, language, and negative features,

$$\mathcal{L}_L = \mathcal{L}_C(\Phi_t^h, e^v, \mathcal{E}^v) + \mathcal{L}_C(\Phi_t^o, e^o, \mathcal{E}^o), \quad (3)$$

where  $e^v$  and  $e^o$  are the target verb and object features, respectively, and  $\mathcal{E}^v$  and  $\mathcal{E}^o$  are the negative verb and negative object feature sets, respectively. More specifically, we rewrite the object term as in Equation (1):  $\mathcal{L}_C(\Phi_t^o, e^o, \mathcal{E}^o) = -(\Phi_t^o)^T e^o + \log \sum_{n=1}^{N_l} \exp((\Phi_t^o)^T \mathcal{E}_n^o)$ , where  $\Phi_t^o$  is the attended object feature that has a similar form as the attended human feature shown in Equation (2),  $e^o$  is the target object feature, and  $N_l$  is the number of negative samples in the negative feature set  $\mathcal{E}^o$ . The human term has a similar form. We show this loss (object term only) in Figure 3 (a). The “Region attention” module outputs a single “Attended human/object feature” for the video frame. This “Attended human/object feature” forms the positive pair with the verb/object phrase in the corresponding language annotation for the frame.

**Self-supervised temporal contrastive loss.** We seek to encourage temporal continuity of the moving objects. We

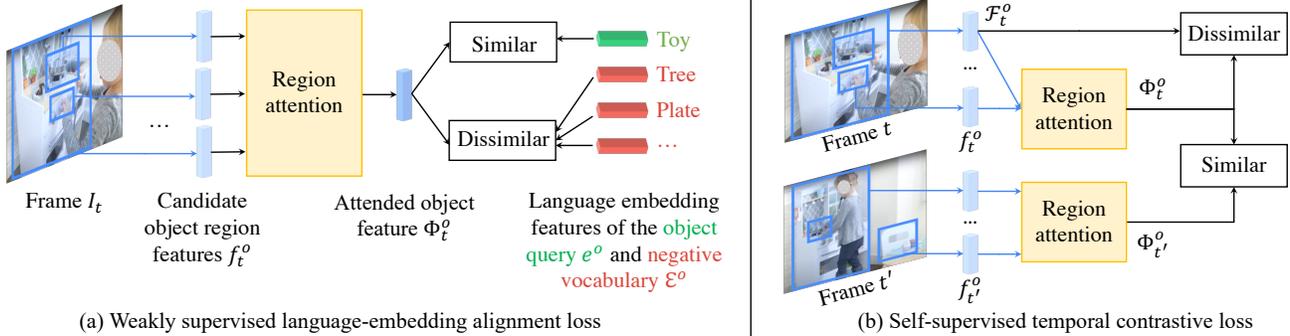


Figure 3: **Weakly supervised contrastive loss.** Our loss jointly aligns features for spatiotemporal regions in a video to (a) a language-embedding feature for an input verb-object query and (b) other spatiotemporal regions likely to contain the target object. This figure only shows object regions. The same mechanism is applied to human regions. (Video credit: KidKraft [23])

also seek to contrast our learned object features against a negative set of visual features corresponding to likely regions for which the target object does not appear. Let  $f_{t'}^o$  be a set of features for another frame from the same video with attention scores  $\hat{\sigma}_{t'}^o$ . We define the temporal contrastive loss  $\mathcal{L}_T$  as the alignment of the attended object feature  $\Phi_t^o$  in a frame  $I_t$  to the target attended object feature  $\Phi_{t'}^o$  in another frame while contrasting against the negative feature set  $\mathcal{F}_t^o$  from frame  $I_t$ . Following the contrastive loss in Equation (1), we define the temporal contrastive loss as:

$$\mathcal{L}_T = \mathcal{L}_C(\Phi_t^o, \Phi_{t'}^o, \mathcal{F}_t^o). \quad (4)$$

Note that the attention scores  $\hat{\sigma}$  here are different from the soft attention scores  $\sigma$  used for the language-embedding alignment loss. In the temporal contrastive loss, we let  $\hat{\sigma}$  be hard attention scores, where only one object region has a score of one while the rest of the regions in the same frame has a score of zero. In practice, we let the object region that has the highest soft attention score have a hard attention score  $\hat{\sigma} = 1$ , which is the most likely target object described in the verb-object query. For the negative feature set  $\mathcal{F}_t^o$ , we randomly select from the remaining object regions in frame  $I_t$  that are not selected by the hard attention. The intuition is that the selected target objects with the highest score from different frames should move consistently through time but should be different from other objects in the same frame. We illustrate this loss in Figure 3 (b).

**Full weakly supervised contrastive loss.** We define the final loss at each frame as the sum of the language-embedding alignment and temporal contrastive losses,

$$\mathcal{L}_{ST} = \mathcal{L}_L + \alpha \mathcal{L}_T, \quad (5)$$

where  $\alpha$  is a hyperparameter. Our loss is minimized when a feature corresponding to a softly selected human/object region  $\Phi_t$  aligns with the language-embedding feature  $e$  and a similar spatiotemporal region  $\Phi_{t'}$  in another frame.

### 3.2. Feature learning

In this section, we briefly introduce the object feature, contextual frame feature, and attended human/object fea-

tures used in Figure 2. See Section A.1 for more details of different types of features.

**Human-guided object feature learning.** To get the features for candidate object regions, we first extract object location proposals in each video frame using Faster R-CNN [39]. We apply ROI pooling over all the layers of the Faster R-CNN feature pyramid network (FPN) to extract the feature descriptors for the object region proposals. Each object region proposal has a feature descriptor  $f_{t,i}^o$  and a bounding box  $b_{t,i}^o$ , as shown in Figure 4.

Human-object interaction is highly correlated with both the human and object features. We assume that the spatial co-occurrence of the human and object regions helps to disambiguate the interacted object. To more effectively encode the human-object interaction, we incorporate knowledge from the human segmentation masks produced by DensePose [1] into the object proposal features. We use ROI pooling to extract a feature  $f_{t,i}^h$  from the human segmentation mask given the object proposal bounding box  $b_{t,i}^o$ . We apply a max-pool operation over the object region features from the FPN feature maps and the human feature maps to obtain the final object proposal feature  $f_{t,i}^o = \max(f_{t,i}^o, f_{t,i}^h)$ .

**Contextual frame feature learning.** Human-object interactions are temporal events and occur over a period of time. To utilize the temporal information from the whole video, we use a soft attention module [51] to learn a contextual feature representation  $x_t$  for each frame. Given a frame feature  $\hat{x}_t$  obtained by passing this frame through a small network, we send  $\hat{x}_t$  to an embedding layer to generate a “query” feature vector  $x_t^{que}$ . For the features of all frames  $\{x_1, \dots, x_T\}$  in the same video, we use two different embedding layers to get “key”  $x_{t'}^{key}$  and “value”  $x_{t'}^{val}$  vectors. We compute the inner product of the “query” and “key” to get a similarity score  $s_{t,t'} = (x_t^{que})^T x_{t'}^{key}$  of the current frame and each frame in the same video. A softmax layer is then applied to the similarity scores to normalize the similarity of each frame to the current frame. The contextual frame feature is obtained by the weighted average over frame “value” features  $x_t = \sum_{t'=1}^T s_{t,t'} x_{t'}^{val}$ .

**Region attended human/object feature learning.** The region attention module computes attention scores for the human/object region proposals to measure their relative relevance to the given verb-object query (Figure 8 in the appendix). For each human region in frame  $I_t$ , we first concatenate its feature representation  $f_{t,i}^h$  with the contextual frame feature  $x_t$  and the verb-object query feature and then pass them through a small network to obtain a score. We apply the softmax function over the scores of all human regions in this frame and get the final human attention scores  $\sigma_t^h$ . Similarly, each object region has an object attention score  $\sigma_t^o$  after applying the softmax function over all object regions. The attention scores are used to aggregate human/object features using Equation (2).

### 3.3. Training objective

In addition to the weakly supervised contrastive loss  $\mathcal{L}_{ST}$ , we propose a sparsity loss  $\mathcal{L}_{spa}$ , and a classification loss  $\mathcal{L}_{cls}$  for weakly supervised learning. Our final training loss for a pair of frames is the sum of all the losses,

$$\mathcal{L}_\theta(t, t') = \mathcal{L}_{ST} + \mathcal{L}_{spa} + \mathcal{L}_{cls}. \quad (6)$$

We describe the sparsity and classification losses next.

**Sparsity loss.** As there are often few humans and objects undergoing the action and object given in the input query, we seek to encourage the attention scores for the human and object proposals each to be high for a single proposal instance and low for all other proposals in each frame. To enable this effect, we introduce a sparsity loss which is defined as the sum of negative log  $L_2$  norms of the human and object attention scores:

$$\mathcal{L}_{spa} = -\log(|\sigma_t^h|_2) - \log(|\sigma_t^o|_2) \quad (7)$$

**Classification loss.** The weakly supervised contrastive loss and sparsity loss enable our model to localize objects and humans given the verb-object query. To make our model retrieve and localize the language input across videos, we add a classification loss to predict whether the current video contains the interaction described in the verb-object query. In the training phase, each video has a ground truth verb-object label and we assign them a label of  $y = 1$ . We randomly select a negative verb-object label from the language features for the entire vocabulary  $\mathcal{E}$  and assign a label of  $y = 0$  to the video and the selected negative verb-object label. The classification loss at frame  $I_t$  is:

$$\mathcal{L}_{cls} = -(y_t \log(p_t^q) + (1 - y_t) \log(1 - p_t^q)), \quad (8)$$

where  $p_t^q = p(y_t | q, x_t)$  is the likelihood of the input video frame  $I_t$  containing the verb-object query  $q$ . Here  $x_t$  is the contextual frame feature of frame  $I_t$ .

### 3.4. Inference

During inference, given a video frame  $I_t$ , we randomly select one verb-object query  $q$  and compute their binary classification score  $p_t^q$  as shown in Equation (8). Since

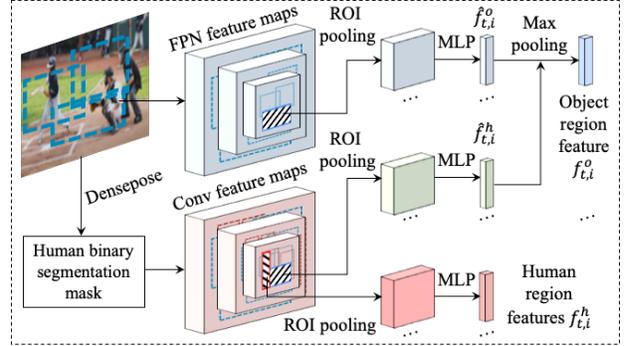


Figure 4: Illustration of extracting human/object features. We learn convolutional filters to encode the Densepose segmentation mask to intermediate features. We obtain the feature of each object region  $f_{t,i}^o$  by combining its ROI pooling features from the FPN feature maps  $\hat{f}_{t,i}^o$  and the human conv feature maps  $\hat{f}_{t,i}^h$ . (Video credit: TheOnDeckCircle [49])

we encourage the matching pairs of video frame and verb-object query to have higher probability during training, score  $p_t^q$  is able to evaluate the probability of the verb-object query appearing in the given frame during inference. We also have an attention score  $\sigma_{t,i}^h$  or  $\sigma_{t,j}^o$  for each human or object region proposal, representing their relevance to the given verb-object query. Thus for each human-object pair, we compute their confidence score as  $c_{t,i,j}^q = p_t^q \times (\sigma_{t,i}^h + \sigma_{t,j}^o) / 2$ . For HOI detection, we predict human and object bounding boxes and their HOI label. For each video frame, we feed in all possible verb-object labels appearing in the dataset and select the verb-object label having the highest confidence score as the HOI label prediction result for each pair of human and object regions.

## 4. Human-object interaction video dataset

Existing human-object interaction datasets either focus on classification [4] or detection in static images [3, 15]. However, human-object interaction is a temporal process and it is more naturally done in video data. Current video datasets, such as Charades [60], EpicKitchens [9], VidVRD [44], VidOR [42], and YouCook [65, 64] are not suitable for human-object interaction detection. First, most of them do not have human bounding box annotations. Second, all objects in a scene are annotated, with annotated objects not necessarily interacting with humans. Furthermore, EpicKitchens and YouCook do not have triplet human-action-object labels. VidVRD and VidOR are for visual relation detection and the relations are not necessarily human-centric. Thus, they cannot be directly used for evaluating video based human-object interaction detection.

Instead, to study the human-object interaction problem in videos, we collect a large, diverse Video dataset of Humans Interacting with Common Objects (V-HICO). Our dataset has a large variety of actions and interacted objects. Our dataset has more videos (6,594) than Epic-Kitchens (432) and YouCook (2,000), with each video containing

Table 1: Evaluation of each component of the proposed model. Phrase (Phr) detection refers to correct localization (0.3 IoU) of the union of human and object bounding boxes while relation (Rel) refers to correct localization (0.3 IoU) of both human and object bounding boxes.

Model	mAP (%)				Recall@1 (%)		Video One Recall@1 (%)		Video All Recall@1 (%)	
	Phr (ko)	Phr (def)	Rel (ko)	Rel (def)	Phr	Rel	Phr	Rel	Phr	Rel
Baseline(add)	40.59	0.45	6.95	0.11	75.98	19.00	90.30	33.22	60.36	7.07
Baseline(cat)	41.86	0.35	11.34	0.11	75.93	19.91	88.49	35.53	61.68	5.92
(cat)+Spa	50.79	1.02	16.23	<b>0.47</b>	79.52	24.24	87.01	38.49	69.74	9.67
(cat)+Spa+Hum	55.60	0.89	15.91	0.29	81.35	22.99	91.61	38.82	70.56	9.55
(cat)+Spa+Hum+Tem	54.42	<b>1.24</b>	16.94	0.30	81.00	25.61	91.12	39.14	69.74	12.68
(cat)+Spa+Hum+Tem+Con	<b>55.90</b>	0.90	<b>18.56</b>	0.26	<b>84.08</b>	<b>30.12</b>	<b>91.94</b>	<b>44.90</b>	<b>75.33</b>	<b>15.95</b>

Table 2: Evaluation of performance on V-HICO compared to methods in [34], [63], [55] and different random baselines. Phrase (Phr) detection refers to correct localization (0.3 IoU) of the union of human and object bounding boxes while relation (Rel) refers to correct localization (0.3 IoU) of both human and object bounding boxes. (ko) and (def) are the known object setting and default setting.

Model	mAP (%)				Recall@1 (%)		Video One Recall@1 (%)		Video All Recall@1 (%)	
	Phr (ko)	Phr (def)	Rel (ko)	Rel (def)	Phr	Rel	Phr	Rel	Phr	Rel
Random	11.24	0.08	0.57	0.00	22.42	4.05	40.79	8.88	6.25	0.49
Random Pretrain	9.58	0.02	0.48	0.00	12.26	3.59	25.33	8.06	1.97	0.33
[34]	32.42	0.14	2.06	0.01	45.75	5.02	71.38	14.14	20.72	0.16
[63]	21.88	0.60	4.83	0.04	55.56	8.04	71.05	16.45	38.49	1.97
[55]	25.34	0.12	4.06	0.05	43.07	5.31	63.16	12.50	24.84	0.49
Ours	<b>55.90</b>	<b>0.90</b>	<b>18.56</b>	<b>0.26</b>	<b>84.08</b>	<b>30.12</b>	<b>91.94</b>	<b>44.90</b>	<b>75.33</b>	<b>15.95</b>

human-object interactions. Furthermore, the new dataset is more challenging with more diverse outdoor scenes compared with Charades, EpicKitchens, and YouCook that either focus on household or kitchen scenes.

Our V-HICO dataset contains 5,297 training videos, 635 validation videos, 608 test videos, and 54 unseen test videos of human-object interactions. To test the performance of models on common human-object interaction classes and generalization to new human-object interaction classes, we provide two test splits, the first one has the same human-object interaction classes in the training split while the second one consists of unseen novel classes. Our training set consists of 193 object classes and 94 action classes. There are 653 action-object pair classes in the training set. The unseen test set contains 51 object classes and 32 action classes with 52 action-object pair classes. All videos are labeled with text annotations of the human action and the associated object. The test set and the unseen test set contain the annotations of both human and object bounding boxes.

Our ‘unseen’ test set (51 unseen object classes) contains 2 classes present in the MSCOCO object vocabulary, 8 present in OpenImages, and 34 present in VisualGenome. We use the object detector pre-trained on MSCOCO, indicating only 2 object classes have been seen during pre-training. Furthermore, our entire dataset has 244 object classes in total. 156 of them are not present in MSCOCO or OpenImages, *e.g.*, ‘javelin’, and hence cannot be detected using detectors pre-trained on those datasets. The object distribution is long-tailed and many objects do not have annotated training data in the publicly available object datasets. Our model provides a way to scale-up to a large

set of objects without relying on bounding box annotations.

## 5. Experiments

We evaluate the ability of our method and baselines to detect human-object interactions on the V-HICO dataset.

### 5.1. Evaluation criteria

We evaluate the proposed method and other approaches under two settings – phrase accuracy and relation accuracy. We denote **phrase** accuracy when the union of the detected human and object bounding boxes matches the union of the ground truth human and object boxes. We denote **relation** accuracy when both the predicted human and object bounding boxes match the ground truth human and object bounding boxes, respectively. Relation accuracy is lower than phrase accuracy since it is more strict on the predicted human and object bounding boxes.

We report the mean average precision (mAP) and Recall in these two setups. For mAP, we follow the settings proposed by HICO-DET [3]. They proposed two different evaluation settings: (1) Known Object setting (**ko**): Given a human-object interaction category, they evaluate the human and object detection only on images containing the target object category. Here we use video frames that contain the target HOI category. (2) Default setting (**def**): Given a HOI category, they evaluate the detection on the full test set. This setting is more challenging because it requires models to distinguish whether an image/frame contains the target HOI category and to localize the target HOI simultaneously. Note that the evaluation metric we used is designed for HOI detection [3], which is a harder problem than lan-

Table 3: Evaluation of our proposed approach, [34], and different random baselines on the unseen test set on V-HICO. The unseen test set consists of 51 classes of objects unseen during training. Evaluation at IoU threshold 0.3.

Model	mAP (%)				Recall@1 (%)		Video One Recall@1 (%)		Video All Recall@1 (%)	
	Phr (ko)	Phr (def)	Rel (ko)	Rel (def)	Phr	Rel	Phr	Rel	Phr	Rel
Random	10.44	0.16	0.74	0.03	14.79	2.11	26.92	5.77	7.69	0.00
Random Pretrain	4.78	0.10	0.42	0.02	14.79	2.82	28.85	5.77	1.92	0.00
Peyre2017 [34]	38.19	0.70	4.79	0.07	43.24	5.41	64.81	12.96	16.67	0.00
Ours	<b>67.21</b>	<b>2.76</b>	<b>25.10</b>	<b>0.66</b>	<b>91.89</b>	<b>31.08</b>	<b>94.44</b>	<b>42.59</b>	<b>85.19</b>	<b>18.52</b>

guage grounding. In language grounding, the query input appears in the video and the models return its corresponding bounding box during test. However, in the **def** setting, the query input does not necessarily appear in the video.

For each frame, we extract the top 10 predicted pairs of human-object bounding boxes based on their score  $c_{t,i,j}^q$  as described in Section 3.4. The predicted human and object bounding boxes are treated as correct if their Intersection-over-Union (IoU) with ground truth human and object bounding boxes is larger than 0.3 for both the phrase and relation accuracy, similar to [34]. We follow HICO-DET [3] and compute the **mAP** over all verb-object classes.

We also report the frame recall of the top-1 prediction. Given a frame and its true verb-object label, we test if the top-1 predicted human-object bounding-box pair matches the ground truth bounding boxes. **Recall@1** is the number of frames where the predictions are correct divided by the number of all frames. We also propose two video recall settings. In **Video One Recall**, if all of the ground truth human-object pairs in one frame are detected, the video is considered correct. Video One Recall is the number of correct videos divided by the number of all videos. In **Video All Recall**, the video is correct only when all of the ground truth human-object pairs in all frames are detected.

## 5.2. Ablation studies on V-HICO

To investigate the effect of each component of our approach, we perform a series of ablation studies on our V-HICO dataset. We report the results in Table 1. We first evaluate our approach when no temporal continuity is enforced in the model during training. To achieve this goal, we omit the temporal contrastive loss  $\mathcal{L}_T$  and the sparsity loss  $\mathcal{L}_{spa}$  during training and do not include the human ROI-pooled feature as part of the object proposal feature. We investigate different ways to merge the human/object region features, verb-object language features, and the frame feature  $\hat{x}_t$  (before the temporal soft attention described in Section 3.2) when computing human/object attention scores. We find that feature addition (**Baseline(add)**) and feature concatenation (**Baseline(cat)**) have similar results.

Next, we evaluate the efficacy of the sparsity loss  $\mathcal{L}_{spa}$ . Without the sparsity loss, empirically we find that the output attention scores are often uniformly distributed across all region proposals. **(cat)+Spa** is the result after using the

sparsity loss based on the feature concatenation baseline; it boosts the performance significantly.

As existing human detectors are quite robust in videos and the spatial location of the human can help localize the interacted object, we evaluate the effect of including the ROI-pooled feature from the human segmentation feature maps to the object region feature (**(cat)+Spa+Hum**). We observe that including the human information when learning object features improves the performance. Next, we evaluate the efficacy of including our self-supervised temporal contrastive loss  $\mathcal{L}_T$ . **(cat)+Spa+Hum+Tem** improves the performance by encouraging the temporal continuity of moving objects. We investigate the effect of the contextual frame feature generated using the soft attention over the entire video. **(cat)+Spa+Hum+Tem+Con** is the result of adding the contextual frame feature and is used as the final result of the proposed model.

To further verify the contribution of human information to object detection results, we add a baseline that localizes the object based on the human spatial prior alone. We first use our model to generate candidate human/object proposals and their confidence scores. We select the human proposal with the highest score as the target human. For each object proposal, we compute its confidence as the inverse distance of its centroid to the human proposal. The mAP on Phr (ko) is 46.67 while ours is 55.90. Note that this baseline’s mAP is not too bad as it uses human/object proposals and human confidence scores from our trained model (**((cat)+Spa+Hum+Tem+Con)**), yet it performs significantly worse than our full model.

## 5.3. Comparison with baselines

Since most existing HOI approaches use supervised learning on static images, we compare with the three most related methods [34, 63, 55] and add two random baselines to compare with our approach in Table 2. “Random” is our model using randomly initialized parameters. “Random Pretrain” is our model with the Faster R-CNN part initialized from the COCO pretrained model and other parts randomly initialized as above.

Since there is no existing weakly supervised human-object interaction detection method for videos, we modify three related weakly supervised baselines using their publicly available code. Peyre *et al.* [34] is a weakly supervised approach for visual relation detection in single



Figure 5: Qualitative predictions of our model with top predicted human bounding boxes (yellow) and object bounding boxes (blue). (Video credit: Dude Chennai [6] and Serious Eats [10])

still images. For fair comparison, our implementation of Peyre *et al.* [34] uses the same human and object bounding boxes and features generated by DensePose and Faster R-CNN, respectively, as our approach. For each human-object bounding box pair, the classifier predicts its probability score of being each human-object interaction class. The human-object bounding box pairs are ranked based on their confidence scores for evaluation. We also compare with Zhou *et al.* [63], a method for video-based object grounding from text. We use the same Faster R-CNN to generate object bounding box proposals. A human detection branch is added using the human proposals generated by DensePose. We further modified a video relation grounding method [55] which also uses the Densepose and Faster R-CNN to generate human and object bounding boxes for a fair comparison.

Table 2 shows the comparisons of our approach and these baseline methods on our V-HICO dataset. Overall, our model outperforms all the baselines as [34] is an image-based method without taking advantage of the video information, [63] optimizes object and human bounding boxes and features separately without explicitly considering human-object interactions, and [55] uses a spatiotemporal region graph that may accumulate errors over time.

#### 5.4. Comparison with baselines on unseen classes

To test the generalization ability of our model on unseen objects, we evaluate our method on 52 unseen verb-object classes from the unseen dataset. Note that there are difficulties when evaluating the Zhou *et al.* [63] and Xiao *et al.* [55] baselines on the unseen dataset as most object labels do not appear in the training set. Zhou *et al.* [63] optimize the word embedding for object and action classes during training; for unseen objects and actions, they do not have an optimized word embedding. While Xiao *et al.* [55] report results on zero-shot relation grounding, they consider the case when the subject-predicate-object triplet is never seen but the separate subject, predicate or object are known during training. However, on our unseen test set, most object labels do not appear in the training set. Thus Xiao *et al.* [55] have the same problem as Zhou *et al.* [63] – they do not have an optimized embedding for unseen words. Thus we only compare our method with Peyre *et al.* [34], “Random”, and “Random Pretrain” on the unseen test set.

Table 3 shows that our method generalises well to new object classes and significantly outperforms the baselines in terms of both the phrase and relation accuracy. Our ap-

proach on the unseen test set is better than the seen test set as the unseen set is smaller and easier (54 videos, most scenes have a single human) than the seen set (608 videos, more challenging scenes with multiple humans or blurry objects). The size of the test set influences some criteria, *e.g.*, mAP is computed over bounding boxes from all videos, thus mAP tends to be lower if there are more videos in the test set.

#### 5.5. Qualitative results

**Human-object interaction detection results.** We present the human and object bounding box predictions of our model in Figure 5. We only show the top 1 human-object bounding box pair. The yellow bounding boxes represent the predicated human bounding boxes while the blue bounding boxes are the predicated object bounding boxes. We find that the proposed weakly supervised method tends to generate large object bounding boxes as learning from weak supervision is challenging. The system must automatically identify and associate video spatiotemporal regions with the provided phrase annotations during training.

**Failure case analysis.** We notice three main failures in our model predictions: (1) when the human prediction is wrong due to incorrect Densepose output (*e.g.*, missed detections when only a small human body part is visible or when multiple people cause heavy occlusions), (2) when the object prediction is incorrect because the object is small, moving, blurry, or partially occluded, and (3) when both detections are incorrect in challenging scenes, *e.g.*, nighttime.

### 6. Conclusion

Weakly supervised HOI detection in videos is a challenging problem, which has not yet received much attention. However, this problem is of great importance as human-object interactions are common in real life with important applications, such as video search and editing, surveillance, and human-robot interaction. In this paper, we introduce a contrastive loss for learning to detect humans and interacted objects in videos given weak supervision. We demonstrate our approach on a new dataset of videos with verb and object phrase annotations. Our approach is a step toward understanding everyday human-object interactions in videos. We hope the proposed dataset and method can facilitate future research in this direction.

**Acknowledgments.** This work was partly supported by the European Regional Development Fund under the project IMPACT (reg. no. CZ.02.1.01/0.0/0.0/15 003/0000468).

# Appendix

In this appendix, we first give the model architecture details and implementation details in Section A. Then we provide the dataset collection details in Section B. In Section C, we show the dataset statistics.

## A. Model architecture details and implementation details

We provide the model architecture details in Section A.1 and implementation details in Section A.2.

### A.1. Feature learning

In Section 3.2 of the main paper, we introduce the object feature, contextual frame feature, and attended human/object features. In this section, we provide more details about the different types of features.

**Verb-object query feature learning.** To extract the feature of the verb-object query described in the main paper (Figure 2), we first map the input verb and object queries to embedded features  $\hat{e}^v$  and  $\hat{e}^o$ , respectively, using the publicly available Google News Word2Vec model [30]. Next, we pass the embedded features through linear mappings  $W_v$  and  $W_o$  to obtain 128-dimensional vectors  $e^v = W_v \hat{e}^v$  and  $e^o = W_o \hat{e}^o$ .

**Human feature learning.** To get the human region features  $f_t^h$  described in the main paper (Figure 2), we first extract candidate human location proposals in each video frame using the publicly available DensePose model [1], which returns a binary segmentation mask of humans in the scene and human bounding-box proposals. As shown in Figure 6, at time  $t$ , each human region proposal  $i$  has a bounding box  $b_{t,i}^h$ . We pass the segmentation mask to a convolutional network to generate human feature maps and then use ROI pooling over the human bounding box  $b_{t,i}^h$  to generate human region features  $f_{t,i}^h$ . The convolutional network consists of a  $7 \times 7$  spatial convolutional layer, followed by ReLU and max-pooling nonlinearities, followed by a  $3 \times 3$  spatial convolutional layer.

**Contextual frame feature learning.** We describe the contextual frame feature learning in the main paper Section 3.2. Here we illustrate the learning process of the contextual frame feature in Figure 7.

Human-object interactions are temporal events and occur over a period of time. To utilize the temporal information from the whole video, we use a soft attention module [51] to learn a contextual feature representation  $x_t$  for each frame. Given a video frame  $I_t$ , we send the frame to Faster R-CNN [39] and extract the final layer of the FasterR-CNN feature pyramid network to obtain an intermediate feature map. We add an average pooling layer after the intermediate feature map and generate a feature vector as the frame

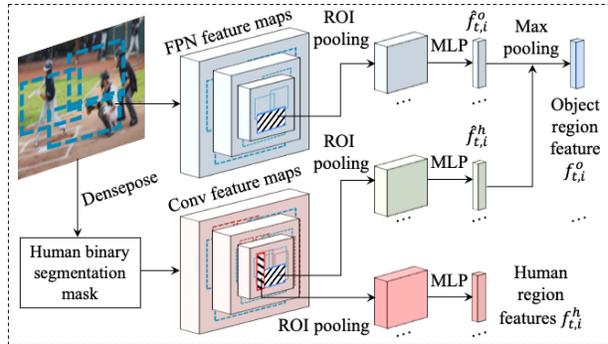


Figure 6: **Illustration of extracting human/object features.** We learn convolutional filters to encode the Densepose segmentation mask to intermediate features. We obtain the feature of each object region  $f_{t,i}^o$  by combining its ROI pooled features from the FPN feature maps  $\hat{f}_{t,i}^o$  and the human conv feature maps  $\hat{f}_{t,i}^h$ . The feature of each human region  $f_{t,i}^h$  is the corresponding ROI pooled feature from the human conv feature maps. (Video credit: TheOnDeckCircle [49])

feature descriptor  $\hat{x}_t$ . Then we send  $\hat{x}_t$  to an embedding layer to generate a “query” feature vector  $x_t^{que}$ . We use the same method to extract the features of all frames in the input video and represent them as  $\{\hat{x}_1, \dots, \hat{x}_T\}$ . For the feature of each frame in the video, we use two different embedding layers to get “key”  $x_{t'}^{key}$  and “value”  $x_{t'}^{val}$  vectors. We compute the inner product of the “query” and “key” to get a similarity score  $s_{t,t'} = (x_t^{que})^T x_{t'}^{key}$  of the current frame and each frame in the same video. A softmax layer is then applied to the similarity scores to normalize the similarity of each frame to the current frame. The contextual frame feature is obtained by the weighted average over frame “value” features  $x_t = \sum_{t'=1}^T s_{t,t'} x_{t'}^{val}$ .

**Region attended human/object feature learning.** To obtain the attended human and object features,  $\Phi_t^h$  and  $\Phi_t^o$ , used in the main paper Figure 2, we first compute an attention score for each human/object region and then aggregate the human/object features based on their attention scores. In Figure 8, we show the details of the region attention module used in the main paper (Figure 2). The region attention module computes attention scores for the human/object region proposals to measure their relative relevance to the given verb-object query. For each human region in frame  $I_t$ , we first concatenate its feature representation  $f_{t,i}^h$  with the contextual frame feature  $x_t$  and the verb-object query feature and then pass them through a small network (consisting of two fully-connected layers with LeakyReLU as the activation function in the middle) to obtain a score. We apply the softmax function over the scores of all human regions in this frame and get the final human attention scores  $\sigma_t^h$ . Similarly, each object region has an object attention score  $\sigma_t^o$  after applying the softmax function over all object regions. The attention scores are used to aggregate human/object features using Equation 2 in the main paper.

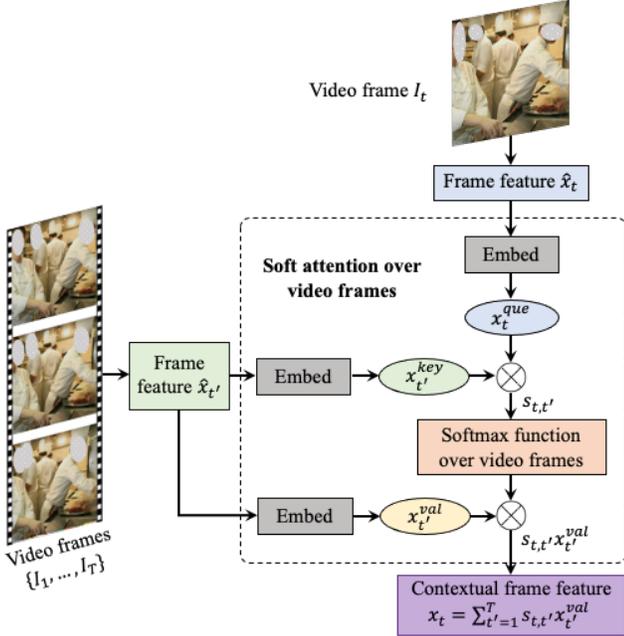


Figure 7: **Illustration of learning contextual frame feature.** Given a frame feature  $\hat{x}_t$  obtained by passing this frame through a neural network, we send  $\hat{x}_t$  to an embedding layer to generate a “query” feature vector  $x_t^{que}$ . For the feature of each frame in the same video, we use two different embedding layers to get “key”  $x_{t'}^{key}$  and “value”  $x_{t'}^{val}$  vectors. We compute the inner product of the “query” and “key” to get a similarity score  $s_{t,t'}$  of the current frame and each frame in the same video. A softmax layer is then applied to the similarity scores to normalize the similarity of each frame to the current frame. The contextual frame feature is obtained by the weighted average over frame “value” features. (Video credit: The Best Gallery Craft [7])

## A.2. Implementation details

Our model is initialized with a ResNext101 Faster R-CNN model, with the RPN pretrained on the COCO dataset from the Detectron library [13]. During training, we select 12 frames from each video and 512 object region proposals (after non-maximum suppression) as object candidate bounding boxes and 25 human bounding boxes for each frame. For the weakly supervised language-embedding alignment  $\mathcal{L}_L$  loss (Equation 3 of the main paper), we compute the loss over 15 sampled negatives from  $\mathcal{E}^v$  for the human term and 15 sampled negatives from  $\mathcal{E}^o$  for the object term, during training.

For the self-supervised temporal contrastive loss  $\mathcal{L}_T$  in each frame  $I_t$ , we compute the loss over 15 sampled negatives from the negative feature set  $\mathcal{F}_t^o$ . In practice, we find that the objects or humans of interest are not always present across all the frames in a video. Some video frames will only show part of the object/human or background. To make the proposed self-supervised temporal contrastive loss more robust to frames that do not contain the mentioned

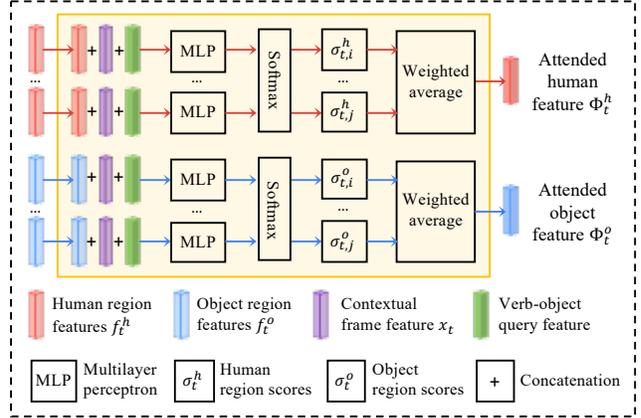


Figure 8: **Illustration of the region attention module.** The region attention module computes attention scores for the human/object region proposals to measure their relative relevance to the given verb-object query. For each human region in frame  $I_t$ , we first concatenate its feature representation  $f_{t,i}^h$  with the contextual frame feature  $x_t$  and the verb-object query feature and then pass them through an MLP to obtain a score. We apply the softmax function over the scores of all human regions in this frame and get the final human attention scores  $\sigma_{t,i}^h$ . Similarly, each object region has an object attention score  $\sigma_{t,j}^o$  after applying the softmax function over all object regions. The attention scores are used to aggregate human/object features as weights in a weighted average given by Equation 2 in the main paper.

human-object interaction, we only use the temporal contrastive losses on 50% frames that have the lowest temporal contrastive losses in each video. The selected frames are more likely to contain the target human and objects.

We used the Adam optimizer [24] with a learning rate of  $1e-4$  and a learning rate of  $1e-6$  for the Faster R-CNN. We use a weight coefficient of  $\alpha = 0.1$  for the temporal contrastive loss  $\mathcal{L}_T$  in Equation 5 of the main paper.

## B. Dataset collection details

We extract the videos from the Moments in Time dataset [31]. The Moments in Time dataset has 800k videos with associated metadata, such as title sentences and tags. Moreover, each video has a manually provided action label, such as “drinking” and “pushing”. We leverage the action labels to help find labels for the human-object interactions from the metadata associated with the videos. We achieve this goal by initially filtering videos to contain the action label in the title sentence or metadata. However, some videos do not have verbs corresponding to human-object interactions, such as “storming” and “erupting”, so we manually discard videos that do not correspond to human actions. We then used the Stanford NLP parser [27] to find videos containing noun phrases after the action label in the title or metadata, and use the resulting noun phrase as the object label. Finally, we remove videos with non-English metadata and manually filter out bad parsing results. After filtering, we

obtained approximately 14,000 videos. We manually filtered out bad examples, such as videos having low frame resolution, wrong language labels, or blurry humans and objects. We finally obtained 6,594 videos in total.

We semi-automatically analyzed the natural language descriptions that accompanied the videos. We do not define a fixed list of HOIs a priori but instead use action-object pairs that appear with a certain frequency in the language captions. By considering more videos with accompanying descriptions, the vocabulary naturally increases.

We collect human and object bounding box annotations using Amazon Mechanical Turk for the test and unseen datasets. We ask each worker to annotate the specific human and object bounding boxes participating in the given human-object interaction label. For each video frame, we collect bounding box annotations from 3 different workers. We average the annotations from each worker to obtain the object bounding box annotations. We assume that there can be multiple people interacting with the given object in a video frame. To obtain the accurate number of humans in the input video frame, we want to cluster the human bounding boxes collected from different workers. The close human bounding boxes are more likely to describe the same human. By counting the number of clusters, we can estimate the number of humans in the input video frame. To do this, we ran an affinity propagation clustering algorithm [32] on all labelled human bounding boxes across multiple workers. We select the clusters which have more than two annotations and average all the annotations within each cluster as the bounding box annotation of that person. We further manually examine the annotated bounding boxes and discard low-quality annotations.

### C. Dataset statistics

Our focus is on video-based, human-centric HOI detection without exhaustively annotating the spatial location of objects in a video at training which is time consuming given the large number of frames in a video. Our dataset consists of 244 different object classes and 99 different action classes. There are 756 verb-object classes in total with diverse human-object interactions.

All the videos are extracted from the Moments in Time dataset [31], which contains short trimmed videos. We semi-automatically analyzed the natural language descriptions that accompanied the videos. By considering more videos with accompanying descriptions, the vocabulary can naturally increase.

We present the dataset statistics in Figure 9, Figure 10, and Figure 11. Figure 9 shows the distribution of objects. We show the top 50 most frequent object classes. Figure 10 shows the distribution of the top 50 most frequent action classes. Figure 11 shows the distribution of the top 50 most frequent verb-object classes.

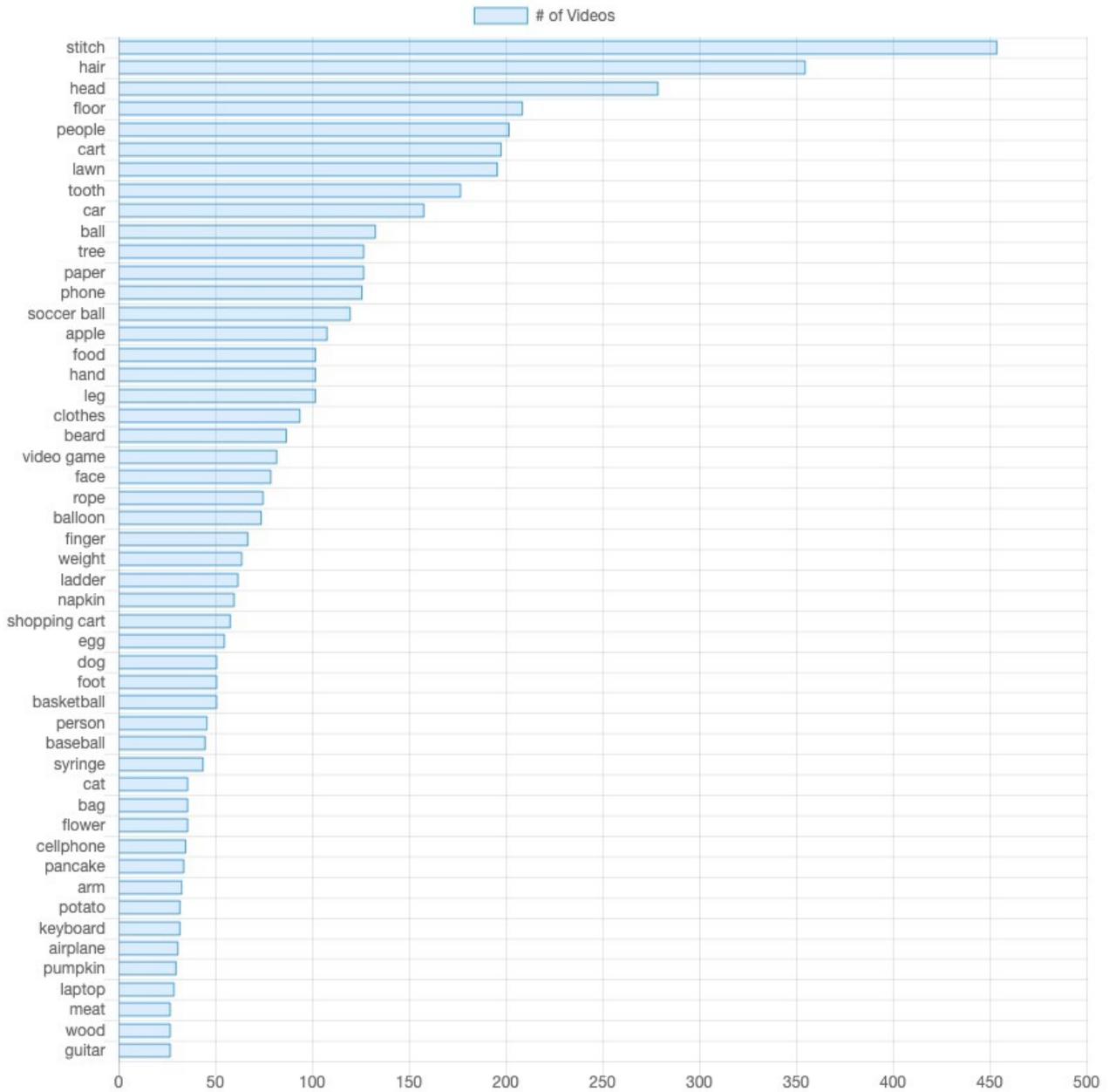


Figure 9: **Distribution of objects in our dataset.** Our dataset consists of 244 different object classes, where for brevity we only show the top 50 in the diagram above.

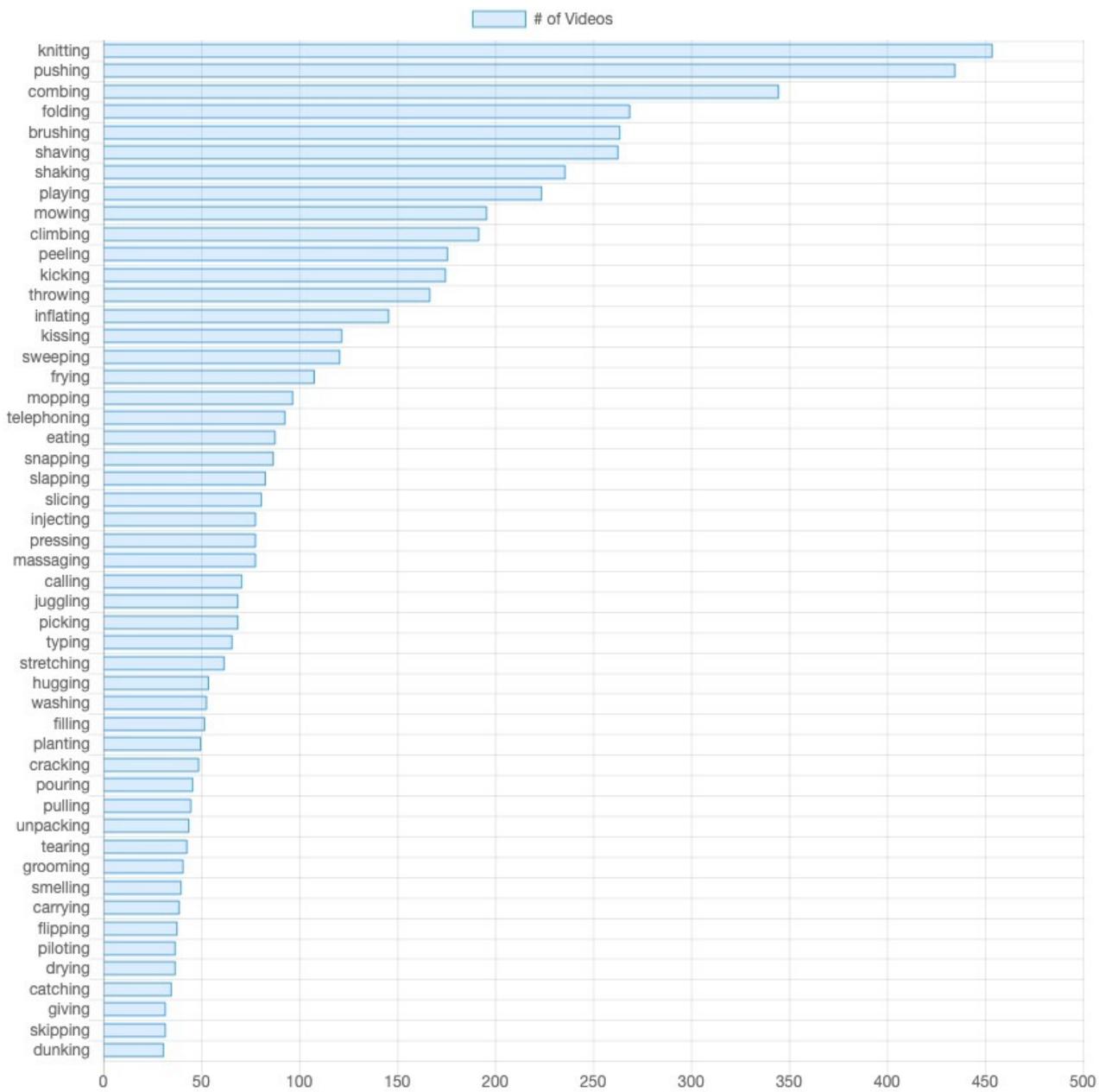


Figure 10: **Distribution of actions in our dataset.** Our dataset consists of 99 different action classes, where for brevity we only show the top 50 in the diagram above.

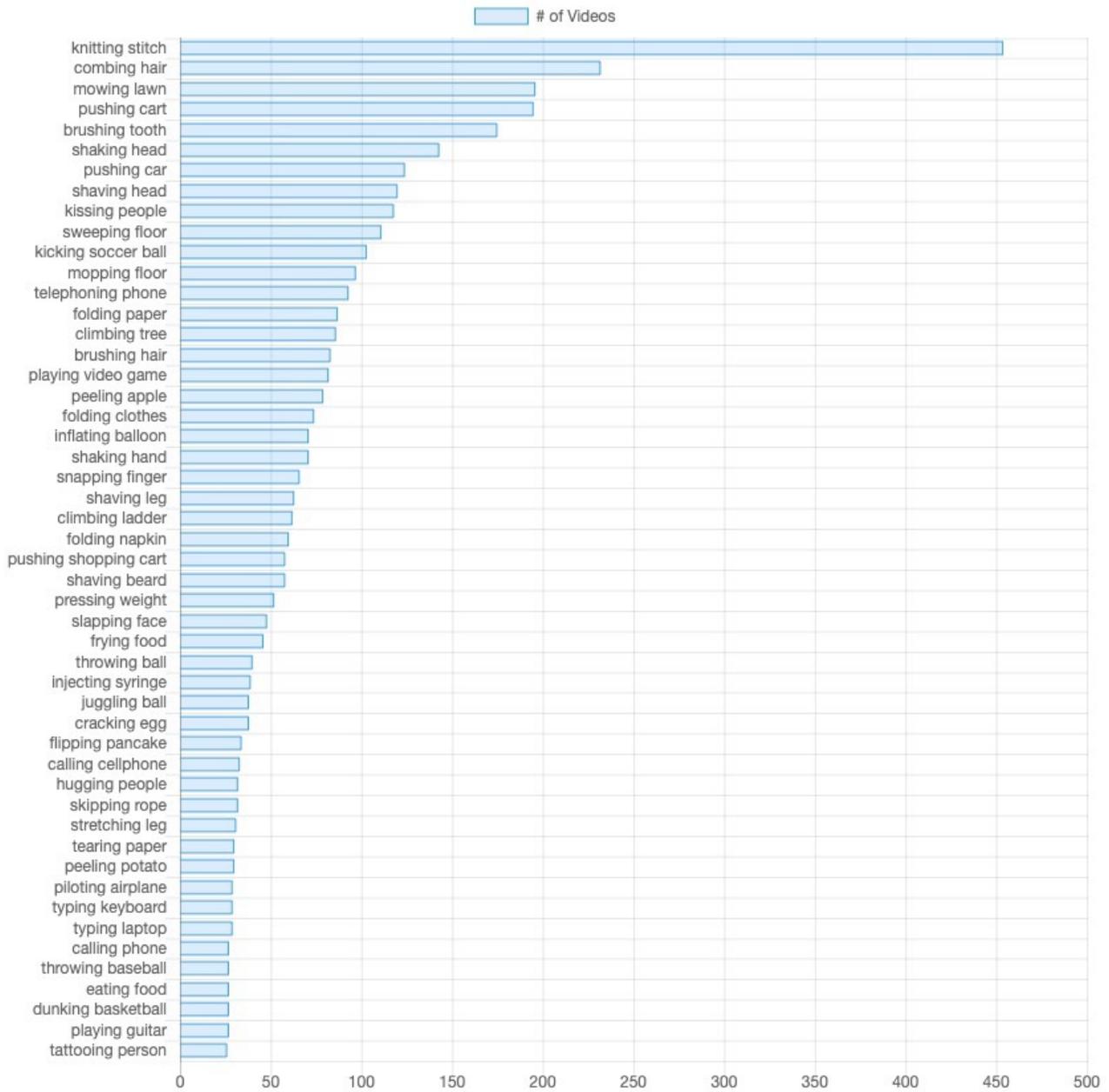


Figure 11: **Distribution of verb-object classes in our dataset.** Our dataset consists of 756 different verb-object classes, where for brevity we only show the top 50 in the diagram above.

## References

- [1] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7297–7306, 2018. 4, 9
- [2] Stephan Baier, Yunpu Ma, and Volker Tresp. Improving visual relationship detection using semantic modeling of scene descriptions. In *International Semantic Web Conference*, pages 53–68. Springer, 2017. 2
- [3] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In *2018 IEEE winter conference on applications of computer vision (wacv)*, pages 381–389. IEEE, 2018. 1, 2, 5, 6, 7
- [4] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiakuan 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, 5
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 3
- [6] Dude Chennai. <https://www.youtube.com/watch?v=kcmAdeypiu>. 2017. 1, 8
- [7] The Best Gallery Craft. <https://www.youtube.com/watch?v=vo07h1vpi54>. 2018. 3, 10
- [8] Bo Dai, Yuqi Zhang, and Dahua Lin. Detecting visual relationships with deep relational networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3076–3086, 2017. 2
- [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 720–736, 2018. 5
- [10] Serious Eats. <https://www.youtube.com/watch?v=hm5cgwiqzu>. 2011. 8
- [11] David F Fouhey, Wei-cheng Kuo, Alexei A Efros, and Jitendra Malik. From lifestyle vlogs to everyday interactions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4991–5000, 2018. 2
- [12] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. TALL: Temporal activity localization via language query. In *ICCV*, 2017. 2
- [13] Ross Girshick, Ilija Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. 10
- [14] 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
- [15] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015. 1, 2, 5
- [16] Yuping Han, Yajing Xu, Shishuo Liu, Sheng Gao, and Si Li. Visual relationship detection based on local feature and context feature. In *2018 International Conference on Network Infrastructure and Digital Content (IC-NIDC)*, pages 420–424. IEEE, 2018. 2
- [17] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019. 3
- [18] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *ICCV*, 2017. 2
- [19] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018. 3
- [20] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. Visual compositional learning for human-object interaction detection. *arXiv preprint arXiv:2007.12407*, 2020. 1, 2
- [21] De-An Huang, Shyamal Buch, Lucio Dery, Animesh Garg, Li Fei-Fei, and Juan Carlos Niebles. Finding ‘it’: Weakly-supervised reference-aware visual grounding in instructional video. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [22] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv preprint arXiv:2102.05918*, 2021. 1
- [23] KidKraft. <https://www.youtube.com/watch?v=caxadtstcng>. 2018. 4
- [24] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 10
- [25] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [26] Kongming Liang, Yuhong Guo, Hong Chang, and Xilin Chen. Visual relationship detection with deep structural ranking. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 2
- [27] Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pages 55–60, 2014. 10
- [28] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9879–9889, 2020. 1, 2
- [29] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by

- Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019. 2
- [30] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv:1301.3781*, 2013. 9
- [31] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Yan Yan, Lisa Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, et al. Moments in time dataset: one million videos for event understanding. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 10, 11
- [32] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. 11
- [33] Julia Peyre, Ivan Laptev, Cordelia Schmid, and Josef Sivic. Detecting unseen visual relations using analogies. In *Proc. ICCV*, 2019. 2
- [34] Julia Peyre, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Weakly-supervised learning of visual relations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5179–5188, 2017. 2, 6, 7, 8
- [35] François Plesse, Alexandru Ginsca, Bertr Delezoide, and Françoise Prêteux. Visual relationship detection based on guided proposals and semantic knowledge distillation. In *2018 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2018. 2
- [36] Bryan A Plummer, Arun Mallya, Christopher M Cervantes, Julia Hockenmaier, and Svetlana Lazebnik. Phrase localization and visual relationship detection with comprehensive image-language cues. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1928–1937, 2017. 2
- [37] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, and Song-Chun Zhu. Learning human-object interactions by graph parsing neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 401–417, 2018. 1, 2
- [38] 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. 1
- [39] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 4, 9
- [40] Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. A dataset for movie description. In *CVPR*, 2015. 2
- [41] Mohammad Amin Sadeghi and Ali Farhadi. Recognition using visual phrases. In *CVPR 2011*, pages 1745–1752. IEEE, 2011. 2
- [42] Xindi Shang, Donglin Di, Junbin Xiao, Yu Cao, Xun Yang, and Tat-Seng Chua. Annotating objects and relations in user-generated videos. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pages 279–287. ACM, 2019. 5
- [43] Xindi Shang, Tongwei Ren, Jingfan Guo, Hanwang Zhang, and Tat-Seng Chua. Video visual relation detection. In *ACM Multimedia*, 2017. 2
- [44] Xindi Shang, Tongwei Ren, Jingfan Guo, Hanwang Zhang, and Tat-Seng Chua. Video visual relation detection. In *ACM International Conference on Multimedia*, Mountain View, CA USA, October 2017. 5
- [45] Liyue Shen, Serena Yeung, Judy Hoffman, Greg Mori, and Li Fei-Fei. Scaling human-object interaction recognition through zero-shot learning. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1568–1576. IEEE, 2018. 2
- [46] Yuhang Song, Wenbo Li, Lei Zhang, Jianwei Yang, Emre Kiciman, Hamid Palangi, Jianfeng Gao, C-C Jay Kuo, and Pengchuan Zhang. Novel human-object interaction detection via adversarial domain generalization. *arXiv preprint arXiv:2005.11406*, 2020. 1
- [47] Yuhang Song, Wenbo Li, Lei Zhang, Jianwei Yang, Emre Kiciman, Hamid Palangi, Jianfeng Gao, and Pengchuan Zhang. Novel human-object interaction detection via adversarial domain generalization. *arXiv preprint arXiv:2005.11406*, 2020. 1
- [48] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelwagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. MovieQA: Understanding Stories in Movies through Question-Answering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [49] TheOnDeckCircle. <https://www.youtube.com/watch?v=kgqltyu6ok>. 2013. 5, 9
- [50] Yao-Hung Hubert Tsai, Santosh Divvala, Louis-Philippe Morency, Ruslan Salakhutdinov, and Ali Farhadi. Video relationship reasoning using gated spatio-temporal energy graph. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017. 4, 9
- [52] Bo Wan, Desen Zhou, Yongfei Liu, Rongjie Li, and Xuming He. Pose-aware multi-level feature network for human object interaction detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9469–9478, 2019. 1, 2
- [53] Suchen Wang, Kim-Hui Yap, Junsong Yuan, and Yap-Peng Tan. Discovering human interactions with novel objects via zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11652–11661, 2020. 1, 2
- [54] Tiancai Wang, Rao Muhammad Anwer, Muhammad Haris Khan, Fahad Shahbaz Khan, Yanwei Pang, Ling Shao, and Jorma Laaksonen. Deep contextual attention for human-object interaction detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5694–5702, 2019. 1, 2
- [55] Junbin Xiao, Xindi Shang, Xun Yang, Sheng Tang, and Tat-Seng Chua. Visual relation grounding in videos. In *European Conference on Computer Vision*, pages 447–464. Springer, 2020. 2, 6, 7, 8

- [56] Bingjie Xu, Yongkang Wong, Junnan Li, Qi Zhao, and Mohan S Kankanhalli. Learning to detect human-object interactions with knowledge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2
- [57] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A large video description dataset for bridging video and language. In *CVPR*, 2016. 2
- [58] Xu Yang, Hanwang Zhang, and Jianfei Cai. Shuffle-then-assemble: Learning object-agnostic visual relationship features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 36–52, 2018. 2
- [59] Ruichi Yu, Ang Li, Vlad I Morariu, and Larry S Davis. Visual relationship detection with internal and external linguistic knowledge distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1974–1982, 2017. 2
- [60] Yuan Yuan, Xiaodan Liang, Xiaolong Wang, Dit-Yan Yeung, and Abhinav Gupta. Temporal dynamic graph lstm for action-driven video object detection. In *ICCV*, 2017. 5
- [61] Yibing Zhan, Jun Yu, Ting Yu, and Dacheng Tao. On exploring undetermined relationships for visual relationship detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5128–5137, 2019. 2
- [62] Hanwang Zhang, Zawlin Kyaw, Jinyang Yu, and Shih-Fu Chang. Ppr-fcn: weakly supervised visual relation detection via parallel pairwise r-fcn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4233–4241, 2017. 2
- [63] Luwei Zhou, Nathan Louis, and Jason J Corso. Weakly-supervised video object grounding from text by loss weighting and object interaction. *arXiv preprint arXiv:1805.02834*, 2018. 6, 7, 8
- [64] Luwei Zhou, Nathan Louis, and Jason J Corso. Weakly-supervised video object grounding from text by loss weighting and object interaction. In *British Machine Vision Conference*, 2018. 5
- [65] Luwei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI Conference on Artificial Intelligence*, pages 7590–7598, 2018. 5
- [66] Penghao Zhou and Mingmin Chi. Relation parsing neural network for human-object interaction detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 843–851, 2019. 1, 2