

Activity Driven Weakly Supervised Object Detection

Zhenheng Yang¹ Dhruv Mahajan² Deepti Ghadiyaram² Ram Nevatia¹
Vignesh Ramanathan²

¹University of Southern California ²Facebook AI

Abstract

Weakly supervised object detection aims at reducing the amount of supervision required to train detection models. Such models are traditionally learned from images/videos labelled only with the object class and not the object bounding box. In our work, we try to leverage not only the object class labels but also the action labels associated with the data. We show that the action depicted in the image/video can provide strong cues about the location of the associated object. We learn a spatial prior for the object dependent on the action (e.g. “ball” is closer to “leg of the person” in “kicking ball”), and incorporate this prior to simultaneously train a joint object detection and action classification model. We conducted experiments on both video datasets and image datasets to evaluate the performance of our weakly supervised object detection model. Our approach outperformed the current state-of-the-art (SOTA) method by more than 6% in mAP on the Charades video dataset.

1. Introduction

Deep learning techniques and development of large datasets have been vital to the success of image and video classification models. One of the main challenges in extending this success to object detection is the difficulty in collecting fully labelled object detection datasets. Unlike classification labels, detection labels (object bounding boxes) are more tedious to annotate. This is even more challenging in the video domain due to the added complexity of annotating along the temporal dimension.

On the other hand, there are a large number of video and image datasets [37, 19, 4, 5, 7] labelled with human actions which are centered around objects. Action labels provide strong cues about the location of the corresponding objects in a scene (Fig. 1) and could act as weak supervision for object detection. In light of this, we investigate the idea of learning object detectors from data labelled only with action classes as shown in Fig. 2.

All images/videos associated with an action contain the object mentioned in the action (e.g. “cup” in the action “drink from cup”). Yuan et al. [50] leveraged this prop-

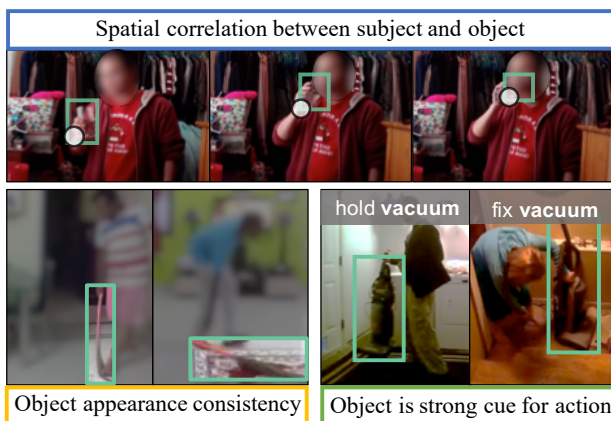


Figure 1: Our framework is built upon three observations we draw: (1) there is spatial dependence between the subject and the interacted object; (2) the object appearance is consistent across different training samples and across different actions involving the same object; (3) the most informative object about the action is the one mentioned in the action.

erty to learn object detection from videos of corresponding actions. However, the actions (“drink from” in above example) themselves are not utilized in this work. On the other hand, the spatial location, appearance and movement of objects in a scene are dependent on the action performed with the object. The key contribution of our work is to leverage this intuition to build better object detection models.

Specifically, we have three observations (see Fig. 1): (1) There is spatial dependence between the position of a person and the object mentioned in the action, e.g. in action “hold cup”, the location of *cup* is tightly correlated with the location of the *hand*. This could provide a strong prior for the object; (2) The object appearance is consistent across images and videos of action classes which involve the object; (3) Detecting the object should help in predicting the action and vice-versa.

The above observations can be used to address one of the main challenges of weakly supervised detection: the presence of a large search space for object bounding boxes during training. Each training image/video has many candidate object bounding boxes (object proposals). In our weakly su-

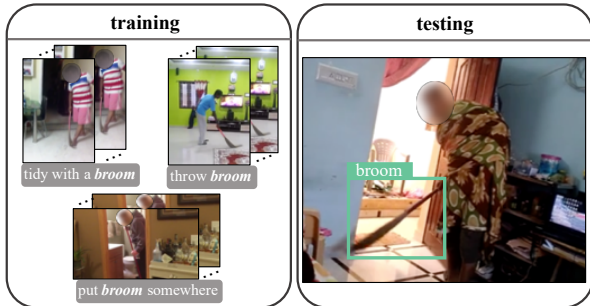


Figure 2: Setting of the action-driven weakly supervised object detection task. Training samples include videos or images with action class labels (left). The inference is conducted on single frame/image for object localization and classification (right).

pervised setting, the only label we have is that, one of these candidates should correspond to the object mentioned in the action. The training algorithm is required to automatically identify the correct object bounding box from this large set of candidates. In our approach, we narrow down this search by incorporating the three observations in our model. In particular, we (1) explicitly learn the spatial prior of objects w.r.t. the human in different actions; (2) train a generic object classifier for modeling object appearances across different actions; (3) jointly learn the action classifier and associated object classifier.

We conducted comprehensive experiments over two video datasets: Charades [37], EPIC KITCHENS [7] and an image dataset: HICO-DET [5]. Our method outperforms the previous methods [3, 50, 41] by a large margin on all datasets. Specifically, we have achieved a 6% mAP boost on Charades compared to current state-of-the-art weakly supervised models for videos. Visualization results and ablation experiments show the effectiveness of each module in our approach.

2. Related Work

In this section, we briefly overview some related research topics and how we are motivated by these works.

Supervised Object detection. Object detection is a very active research topic in the computer vision field. There has been significant progress in the recent years with the advances of deep learning. R-CNN [15] is the first work that introduces CNN features to object detection. A sequence of later works are developed based on R-CNN. Fast R-CNN [14] accelerates R-CNN by introducing an ROI pooling layer and improve the performance by applying proposal classification and bounding box regression jointly. Faster R-CNN [33] further improves the speed and accuracy by replacing the proposal generation stage with a learnable network: region proposal network and the whole framework is trained in an end-to-end fashion. Mask R-CNN [21] proposed to add a segmentation branch and achieved the state-of-the-art (SoTA) performance. All the methods require full object bounding box annotations and mask R-CNN requires

dense segmentation labels.

Weakly supervised object detection. The fully supervised object detection methods rely heavily on large scale bounding box annotations, which is inefficient and labor consuming. To alleviate this issue, there have been various weakly-supervised works [6, 39, 3, 26, 24, 32, 36, 47, 30, 2, 38, 55, 51, 40, 54, 53, 10, 35, 43, 52] that leverage the more efficient image-level object class annotations. Weakly supervised deep detection networks (WSDDN) [3] proposed an end-to-end architecture to perform region selection and classification simultaneously. It is achieved by separately performing classification and detection headers and the supervision comes from a combination classification score. ContextLocNet [26] further improves WSDDN by taking contextual region into consideration. Beyond the image domain, another line of research works [28, 45] try to leverage the temporal information in videos to facilitate the weakly supervised object detection. Kwak *et al.* [28] proposed to discover the object appearance presentation across videos and then track the object in temporal space for supervision. Wang *et al.* [45] perform unsupervised tracking on videos and then cluster similar deep features to form visual representation.

Yuan *et al.* [50] proposed a much more efficient action-driven weakly supervised object detection setting which aims to learn the object appearance representation given only videos with clip-level action class labels. They proposed to first extract spatial features from object proposals. The features are then updated using long short-term memory (LSTM) [22] applied on neighboring frames. The frame-level object classification loss is computed on the updated features. We implemented the same setting as in [50]: pipeline trained on videos/images with only action labels and test on images. Unlike TD-LSTM [50] that only leverages object class information, we propose to jointly exploit both action and object class labels. Considering all actions are interactions between person and objects, we incorporate human pose estimation into the framework.

Activity recognition There are a variety of works in the field of action recognition [44, 29, 17, 1, 12, 11, 13]. Maji *et al.* [29] train action specific poselets that are then classified using SVMs. The contextual cues are captured by explicitly detecting objects and exploiting action labels of other people in the image. R^{*}CNN [17] proposed to implicitly model the main objects. The features from both person region and object proposal regions are extracted and a fusion of classification scores from these two types of features is used for action classification loss. R^{*}CNN showed that the most informative object in the scene is the object mentioned in action class. We are inspired by the similar idea to jointly consider the action and object labels.

Human-object interaction There are mainly two tasks in the human-object interaction (HOI) topic: HOI recogni-

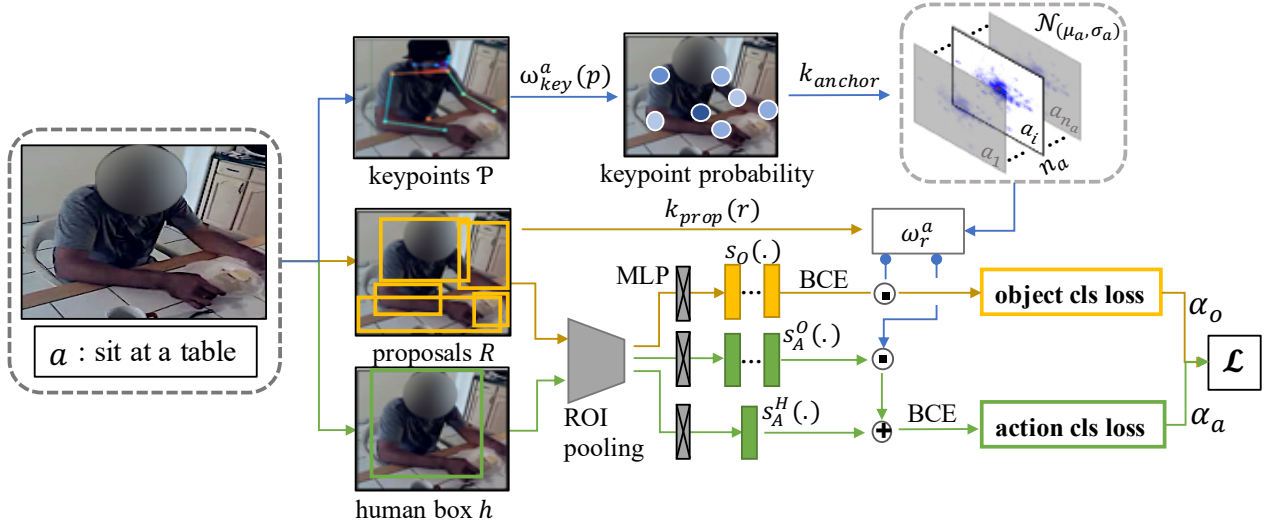


Figure 3: The diagram of our framework. There are three streams in the proposed framework: object spatial prior module (colored in blue), object classification stream (colored in yellow) and action classification module (colored in green). We incorporate human keypoint detection into the framework and jointly leverage action and object labels.

dition and HOI detection. HOI recognition aims at recognizing the interaction between subject and object. There have been a surge of works [8, 20, 49, 34] on HOI recognition since 2009. HOI detection task aims at localizing subject and object and also recognizing HOI class. Chao *et al.* [5] proposed a three-stream network for this task, one stream each for person detection, object detection and person-object pair wise classification respectively. Gkioxari *et al.* [16] model the interaction with shared weights between human centric branch and interaction branch. Kalogeton *et al.* [25] proposed to jointly learn the object and action (e.g. *dog running*). All works have shown that jointly learning the object, person localization and HOI/action classification benefits the performance.

3. Approach

The main challenge of weakly supervised detection is the lack of bounding box information during training and the availability of only image/video level labels. This problem is typically handled in a Multiple Instance Learning (MIL) setting [3, 6], where the training method implicitly chooses the best bounding box from a set of candidate proposals in the image/video to explain the overall image/video label.

However, in practice the number of candidate object proposals can be quite large, making the problem challenging. In our work, we address this issue by imposing additional constraints on the choice of the best object bounding box based on the location prior of the object w.r.t. the human and the importance of the chosen object proposal for action classification. In practice, we model each of these as three different streams in our model which finally contribute to a single action classification loss and an object classification loss. Note that in our work we assume that a pre-

trained person detection model and human keypoint detection model are available to extract the signals needed for capturing human-object dependence.

3.1. Framework

Formally, for a training sample (video clip or image), the action label a is provided. The action a belongs to a pre-defined set of actions $a \in \mathcal{A}$, which is of size n_a : $|\mathcal{A}| = n_a$. We assume that all human actions are interactive and there is one object involved in each action. For example, the object *cup* appears in the action *holding a cup*. The object class associated with action a is denoted by o_a and there are n_o object classes in total: $o_a \in \mathcal{O}$, $|\mathcal{O}| = n_o$.

A pre-trained human detector [21] and pose estimation network [46] are used to extract human bounding box h and keypoint locations $k(p), p \in \mathcal{P}$ where \mathcal{P} represents the set of human keypoints. For training samples with multiple people, we pick the detection result with the highest detection confidence. The object proposals R are extracted. We remove proposals with high overlap ($IoU > \theta_h$) with human region h and we keep the top n_r with highest confidence.

Our model has three streams which are explained in detail in the next sections. An overview of our models is shown in Fig. 3. The first stream models the spatial prior of the object w.r.t. to human keypoints in each action. The prior is used to construct an object classification stream which weights the object classification losses of different proposals in an image/video. The weights and features from the object proposals, along with features of human bounding box, are used to construct an action classification loss. The combined loss from action classification and object classification is minimized during training.

3.2. Object spatial prior

The object spatial prior is modeled in two stages: (1) given an action class a and keypoint detection results \mathcal{P} , we estimate an anchor location based on a weighted combination of the keypoint locations; (2) given the action class and the anchor position, the position of the object is modeled as a normal distribution w.r.t. the anchor point. This is based on our observation that for a given action, certain human keypoints provide strong location priors for the object locations (“hand” for drinking from a cup, “foot” for kicking a ball etc.).

The anchor location k_{anchor} is calculated as a weighted sum of all keypoint locations. The keypoint weight is modeled with a probability vector $w_{key}^a(p)$, $p \in \mathcal{P}$ for the action class a .

$$k_{anchor} = \sum_{p \in \mathcal{P}} w_{key}^a(p)k(p) \quad (1)$$

where $k(p)$ is the detected position of the keypoint p in the training image/video. Given the action class a , the weight of object location w.r.t. the anchor location is modeled with a learned normal distribution: $\mathcal{N}_{(\mu_a, \sigma_a)}$ $\mu_a \in \mathbb{R}^2, \sigma_a \in \mathbb{R}^{(2 \times 2)}$. μ_a represents the mean location of the object w.r.t. the anchor and σ_a represents the variance. This distribution is used to calculate the object location probabilities of different locations. Specifically, the probability of an object being at the location of a proposal $r \in R$ for an action class a is

$$w_r^a = \mathcal{N}_{(\mu_a, \sigma_a)}(k_{prop}(r) - k_{anchor}) \quad (2)$$

where $k_{prop}(r)$ is the center of the proposal r . Note that the distributions w_{key} , $\mathcal{N}_{(\mu_a, \sigma_a)}$ are learned automatically during training.

3.3. Object classification

For each proposal $r \in R$ in a training sample, we compute an object classification score for each object o : $s_O(r; o)$. Here s_O corresponds to an ROI-pooling layer followed by a Multi Layer Perceptron (MLP) which classifies the input region into n_o object classes. Apart from only leveraging image-level object labels for classification [3, 26], the spatial location weights from previous section are also used to guide the selection of the object proposal. Formally, the binary cross-entropy (BCE) loss is calculated on each proposal region, against the image-level object class ground truth. The BCE losses are weighted by the location probabilities of different proposals and the weighted sum is used to compute object classification loss:

$$\begin{aligned} \mathcal{L}_{obj} &= -\frac{1}{n_r} \sum_{r \in R} w_r^a \cdot \mathcal{L}_o(r), \\ \mathcal{L}_o(r) &= \frac{1}{n_o} \sum_{o \in \mathcal{O}} y_o \log(P(o|r)) + (1 - y_o) \log(1 - P(o|r)), \\ P(o|r) &= \frac{\exp(s_O(r; o))}{\sum_{o \in \mathcal{O}} \exp(s_O(r; o))}, \end{aligned} \quad (3)$$

where y_o is the binary object classification label for the object o . Note that y_o is non-zero only for the object mentioned in the action corresponding to the image/video.

3.4. Action classification

For the task of action recognition, especially for interactive actions as in our task, both the person and the object appearances are vital cues. As indicated in [17], the spatial location of the most informative object can be mined from action recognition task. We incorporate a similar idea into the action classification stream by fusing features from the proposal regions and person region. Formally, for a training instance with action label a , the appearance features of both person region h and proposal regions R are extracted, and then classified to n_a -dimension action classification scores: $s_A^O(r; a)$, $r \in R$ and $s_A^H(h; a)$. Here s_A^H , s_A^O correspond to an ROI-pooling layer followed by a Multi Layer Perceptron (MLP). The weights and biases of the MLP are learned during training. The final proposal score is computed as an average of action classification scores weighted by the spatial prior probabilities as in the previous section. This ensures that only scores from the most relevant proposals are given a higher weight. The sum of action classification scores from object proposals and person regions is used to compute the final BCE action classification loss. The loss is computed as follows:

$$\begin{aligned} \mathcal{L}_{act} &= -\frac{1}{n_a} \sum_{a \in \mathcal{A}} y_a \log(P(a)) + (1 - y_a) \log(1 - P(a)), \\ P(a) &= \frac{\exp(s_A^H(h; a) + \sum_{r \in R} w_r^a s_A^O(r; a))}{\sum_{a \in \mathcal{A}} \exp(s_A^H(h; a) + \sum_{r \in R} w_r^a s_A^O(r; a))}, \end{aligned} \quad (4)$$

where y_a is the binary action classification label for the action a .

3.5. Temporal pooling for videos

Our experiments are conducted on both video and image datasets, thus the training samples can be video sequences or static images with action labels. For models trained with video clips, we adopt a few pre-processing steps and also pool scores across the temporal dimension to improve person detection and object proposal quality. Formally, n frames are uniformly sampled from the training clip, followed by person detection and object proposal generation for the sampled frames. The object proposals as well as person bounding boxes across the frames are then connected by an optimization based linking method [18, 48] to form object proposal tubelets and person tubelets respectively. We observed that temporal linking of proposals avoids spurious proposals and leads to more robust features from the proposals. These are fed as inputs into the object classification and action classification streams. Temporal pooling is used to aggregate classification scores across the person and object tubelets. The pooled scores are finally used for loss computation as before.

3.6. Loss terms

The combined loss is a weighted sum of both classification loss terms.

$$\mathcal{L} = \alpha_o \mathcal{L}_{obj} + \alpha_a \mathcal{L}_{act} \quad (5)$$

The hyper-parameters α_o and α_a are weights to trade off the relative importance of object classification and action classification in the pipeline.

3.7. Inference

During testing, object proposals are firstly extracted on the test sample. The trained object classifier (s_O) is applied on each proposal region to obtain the object classification scores ($P(o|r)$). Then the non-maximal suppression (NMS) is applied and the object proposals with higher classification scores than the threshold are preserved as detection results.

4. Experiments

Our method is applicable to both video and image domains. We require only human action label annotations for training. Object bounding box annotations are used only during evaluation. Code will be released in the Github repository ¹.

Video datasets: The Charades dataset [37] includes 9,848 videos of 157 action classes, among which, 66 are in-teractive actions with objects. There are on average 6.8 action labels for a video. The official Charades dataset doesn't provide object bounding box annotations and we use the annotations released by [50]. In the released annotations, 1,812 test videos are down-sampled to 1 frame per second (fps) and 17 object classes are labeled with bounding boxes on these frames. There are 3.4 bounding box annotations per frame on average. We follow the same practice as in [50]: train on 7,986 videos (54,000 clips) and evaluate on 5,000 randomly selected test frames from 200 test videos.

The EPIC-KITCHENS [7] is an ego-centric video dataset which is captured by head-mounted camera in different kitchen scenes. In the training data, the action class is annotated for 28,473 trimmed video clips and the object bounding boxes are labeled for 331 object classes. As the object bounding box annotations are not provided for the test splits, we divide the training data into training, validation and test parts. The top 15 frequent object classes (which are present in 85 action classes) are selected for experiments, resulting in 8,520 training, 1,000 validation and 200 test video clips. We randomly sample three times from each training clip and generate 28,560 training samples. We also randomly sample 1,200 test frames from the test clips.

Image dataset The HICO-DET dataset [5] is designed for human-object interaction (HOI) detection task. This dataset includes 38,118 training images and 9,658 test images. The human bounding box, object bounding box and

an HOI class label are annotated for both training and test images. In total, there are 80 object classes (e.g. *cup*, *dog*, etc.) and 600 HOI classes (e.g. *hold cup*, *feed dog*, etc.). We filter out all samples with "no_interaction" HOI labels, interaction class with less than 20 training samples and all "person" as object class samples. This results in 32,100 training samples of 510 interaction classes and 79 object classes. We use the HOI labels as action class labels during training and the object bounding box annotations are used only for evaluation. Unlike Charades where the interactions mostly happen between one subject and one object, there are cases where multiple people interact with one object (e.g. "boarding the airplane") and one person interacts with multiple objects (e.g. "herding cows"), which makes it more challenging to learn the object appearance.

We report per-class *average precision* (AP) at *intersection-over-union* (IoU) of 0.5 between detection and ground truth boxes, and also mean AP (mAP) as a combined metric, following the tradition of [50]. We also report *CorLoc* [9], a commonly-used weakly supervised detection metric. CorLoc represents the percentage of images where at least one instance of the target object class is correctly detected (IoU>0.5) over all images that contain at least one instance of the class.

4.1. Implementation details

We use VGG-16 and ResNet-101 pre-trained on ImageNet dataset as our backbone feature extraction networks. All *conv* layers in the network are followed with *ReLU* activation except for the top classification layer. Batch normalization [23] is applied after all convolutional layers. In order to compute the classification scores (s_O, s_A^H, s_A^O), three branches are built on top of the last convolutional block. Each branch consists of ROI-pooling layer and 2-layer multiple layer perception (MLP) of intermediate dimension of 4096. The threshold for removing person proposal regions is set as $\theta_h = 0.5$. Selective search [42] is used to extract object proposals for all our experiments.

The Adam optimizer [27] is applied with learning rate of 2×10^{-5} and batch size of 4. The loss weights are set as $\alpha_a = 1.0$, $\alpha_o = 2.0$. The number of sampled frames in a clip is set as $n = 8$ and the number of proposals is set as $n_r = 700$. The whole framework is implemented with PyTorch [31]. We train on a single Nvidia Tesla M40 GPU. The whole training converges in 20 hrs. More details of implementation are presented in supplemental material.

4.2. Influence of modeling spatial location of object

Unlike many existing methods for weakly supervised object detection, our framework explicitly models the spatial location of the object w.r.t. to the detected person and encodes it into two different loss functions in Eq. 3, 4. We explore the effect of modeling this spatial prior through different distributions and its contribution to each of the loss terms.

¹<https://github.com/zhenheny/Activity-Driven-Weakly-Supervised-Object-Detection>

Table 1: Detection performance of different variants on Charades

Spatial prior	Loss term	mAP	CorLoc
Center	action+object	3.43	34.27
Grid	action+object	4.32	36.94
Normal (μ)	action+object	6.27	42.36
Normal (σ)	action+object	4.86	38.05
Normal ($\mu + \sigma$)	action	2.61	31.60
Normal ($\mu + \sigma$)	object	5.86	39.24
Normal ($\mu + \sigma$)	action+object	8.76	47.91

The different distributions include: (a) normal distribution, (b) a fixed grid of probability values, where we make a discrete version of spatial prior module by pre-defining a 3×3 grid around the keypoint, and (c) a simple center prior where we penalize object detections farther away from the center of the object. Note that, we totally removed person detection bounding box and pose estimation in the center prior baseline. For this baseline, we use the frame center as the anchor location \mathcal{L} and learn the μ_a and σ_a .

We also experimented with learning distribution mean only (μ_a), learning variance only (σ_a) and joint learning of mean and variance ($\mu_a + \sigma$) for the normal distribution. We also experimented with using only object classification or action classification loss.

The quantitative results of VGG-16 as the backbone network are presented in Tab. 1 for different ablation settings. First, we observe that a learnable grid-based or normal distribution for the anchor location outperforms a simple heuristic choice of the image center as the anchor. We also see that the normal distribution, where both mean and variance are learned for each action-object pair leads to better results compared to the other settings. This shows that good modeling of the object spatial prior w.r.t. human in an action provides strong cues for detection. We also notice that jointly modelling both action and object classification achieves the best result.

We also visualize the learned distribution of object location probabilities from the prior module for a few sample videos/images in Fig. 4. The learned distribution often has large probability weights around the object mentioned in the action. For example, in the first two columns of the visualization, it is much easier to localize the object with the cues from the heatmap. However, we also note that this distribution is less useful for actions where there is no consistent physical interaction between the human and the object. This is shown in the last column of the figure, for actions like “watching television” and “flying kite”. Our approach reports relatively low mAP performance on such object classes (Tab. 2 and Tab. 3).

4.3. Comparison with existing methods

We compare our method with other weakly supervised methods and their variants: (1) WSDDN [3]; (2) Context-LocNet [26]; (3) PCL [41]; action-driven weakly supervised object detection method: (4) TD-LSTM [50] and (5) R*CNN [17] which is designed for action recognition with

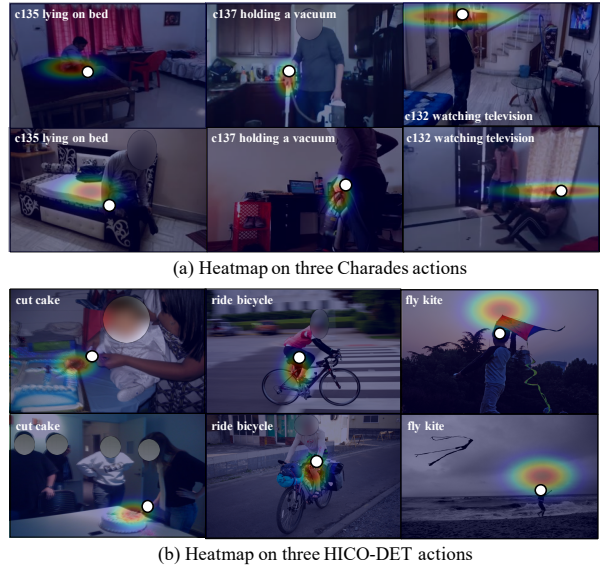


Figure 4: Visualization of learned object location probability w.r.t. selected person keypoint. The heatmap represents the probability of object location (brighter color represents larger probability value) and the white circle represents the selected keypoint.

awareness of the main object. We used the main object bounding box as the object detection result. R*CNN is pre-trained on Pascal-action dataset and then finetuned on Charades or HICO-DET dataset. Note that existing methods (1), (2), (3), (4) do not use person bounding box or keypoint detection results in their model unlike our method. While (5) uses person bounding box, it doesn’t use person keypoints.

The person detection and pose models used in our model were trained only once and kept fixed during training. The annotations required to train the person models are very inexpensive in comparison to fully supervised models which need bounding box annotation for every object class. The resource demands of annotating person bounding boxes and pose is amortized across all object classes. However, since these models are not used in traditional weakly supervised methods, we enable fair comparison by constructing variants of PCL and R*CNN: (6) R*CNN with spatial prior and (7) PCL with spatial prior, where we replace the max pooling in R*CNN and mean pooling in PCL with a weighted sum where the weights are computed from spatial prior distribution as in our implementation (more details are presented in supplemental material).

Results from TD-LSTM [50] are shown only for Charades, since it is a video-specific model and code is not available. Also, we report results from weakly-supervised models whose code is available or whose results on Charades, HICO-DET or EPIC KITCHENS datasets is readily available. Also, note that many methods such as [10, 47, 2] are built on top of the vanilla WSDDN method by adding signals such as segmentation, contextual information, in-

Table 2: AP performance (%) on each object class and mAP (%) comparison with different weakly supervised methods on Charades.

Methods	bed	broom	chair	cup	dish	door	laptop	mirror	pillow	refri	shelf	sofa	table	tv	towel	vacuum	window	mAP(%)
WSDDN [3]	2.38	0.04	1.17	0.03	0.13	0.31	2.81	0.28	0.02	0.12	0.03	0.41	1.74	1.18	0.07	0.08	0.22	0.65
R*CNN [17]	2.17	0.44	2.03	0.31	0.08	0.77	2.64	0.32	1.24	2.36	0.82	1.41	0.65	0.72	0.07	0.65	0.17	0.99
ContextLocNet [26]	7.40	0.03	0.55	0.02	0.01	0.17	1.11	0.66	0.00	0.07	1.75	4.12	0.63	0.99	0.03	0.75	0.78	1.12
TD-LSTM [50]	9.19	0.04	4.18	0.49	0.11	1.17	2.91	0.30	0.08	0.29	3.21	5.86	3.35	1.27	0.09	0.60	0.47	1.98
PCL [41]	4.62	1.07	2.21	1.26	1.08	2.49	3.61	5.13	1.34	4.46	3.29	5.61	3.84	3.26	1.17	1.43	2.27	2.83
R*CNN + prior	6.82	3.64	5.39	3.25	2.47	3.36	5.27	1.07	2.38	6.34	3.29	5.72	4.09	1.03	1.26	3.41	0.86	3.50
PCL + prior	10.57	5.63	8.24	3.52	3.71	5.63	6.86	4.96	5.23	11.39	4.88	10.46	6.32	3.53	4.06	4.89	3.07	6.05
Ours-vgg-16 (w/o prior)	6.71	2.32	5.48	2.49	1.04	3.60	4.02	3.42	4.39	7.76	3.15	7.43	3.26	1.62	0.89	2.24	1.23	3.60
Ours-vgg-16	14.92	10.23	13.08	7.65	5.21	6.44	8.65	4.79	9.14	18.07	7.29	17.21	8.46	2.37	5.46	7.23	2.64	8.76
Ours-ResNet-101	16.54	11.63	14.87	8.62	6.73	8.29	11.32	4.96	9.81	19.24	9.03	18.49	9.86	3.05	6.48	8.08	3.02	10.03

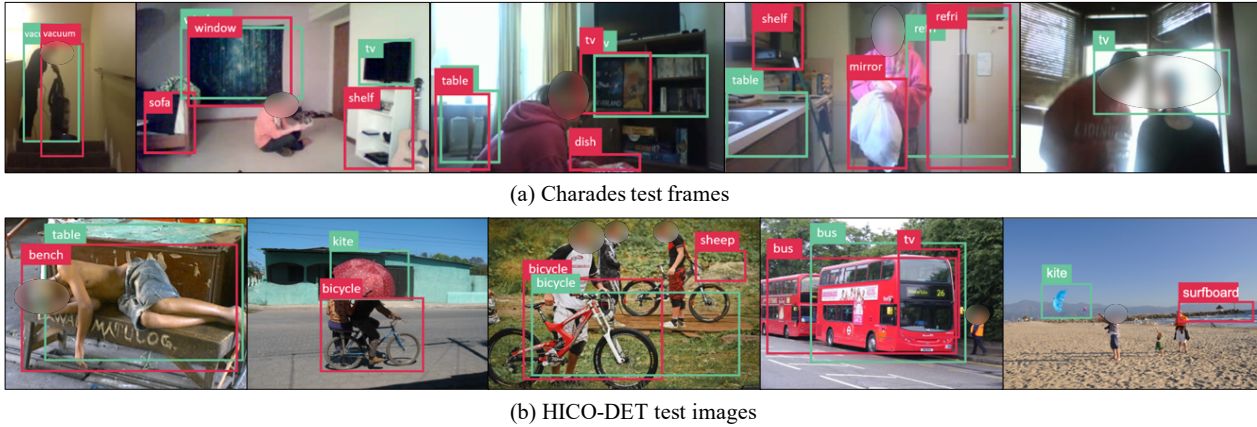


Figure 5: Qualitative detection results on (a) Charades test frames and (b) HICO-DET test images. Red bounding boxes denote our results and green bounding boxes denote results of PCL [41]

stance refinement, *etc.* and these ideas are complementary to the ones presented in this work, and can be added to our model to achieve better results. The per-class AP and combined mAP performances on the two datasets are presented in Tab. 2 and Tab. 3 respectively. 10 object classes on HICO-DET are randomly selected and presented.

On Charades dataset, our method achieves 6% mAP boost compared to PCL [41]. Our method performs better on object classes like broom, refrigerator, vacuum, *etc.*. The spatial prior patterns of the interactions involving these object classes are more predictable and thus the prior modeling benefits our approach more than on other object classes. For object like tv, the spatial prior pattern of the interaction (*e.g. watch tv*) is more diverse and thus difficult to model, resulting only a small boost in mAP. The same performance pattern also applies to HICO-DET dataset. On the object class kite, our method slightly performs inferior to the baseline method.

We observe that the spatial prior from our model is effective in localizing the object during training even when combined with other models such as R*CNN [17] and PCL [41]. R*CNN with spatial prior modeling outperforms TD-LSTM, which is specifically designed for the action driven weakly supervised object detection task.

We also report our model’s performance without the spatial prior module (ours (w/o prior)). This variant of the

model doesn’t require any person bounding box and key-point information, and is directly comparable to existing weakly supervised methods. We note that even without these signals, our model can outperform existing methods. This can be attributed to the fact that our model jointly uses both action and object labels during training. It identifies the object bounding box which can both help action classification and object classification during training.

The qualitative comparison between our method and PCL is presented in Fig. 5. Our approach localizes the object more accurately. Multiple object classes and multiple instances can be detected through our trained object classification stream. The last column shows our failure cases. On Charades, both PCL [41] and our method fails to detect the windows and on HICO-DET, our method fails to localize the kite. One possible reason is that actions like “watch out of the window” do not have direct human-object interaction.

Our approach is also extended to ego-centric EPIC KITCHENS dataset. Since human keypoints are not visible in this dataset, we applied “center” spatial prior modeling used in Sec. 4.2. As the camera is fixed with respect to the human, the anchor location is already implicitly modeled by this center prior. We compare with R*CNN [17] and PCL [41] on the 1,200 test frames. Egocentric videos have a strong prior for object spatial locations and hence our method is able to outperform other methods in Tab. 4.

Table 3: AP performance (%) on selected object classes and mAP (%) comparison with other weakly supervised methods on HICO-DET.

Methods	apple	bicycle	bottle	chair	cellphone	frisbee	kite	surfboard	train	umbrella	mAP(%)
R*CNN [17]	1.13	3.26	1.57	2.35	1.47	1.02	0.32	2.70	2.86	3.04	2.15
WSDN [3]	1.46	5.19	1.52	3.87	2.02	2.44	1.15	2.86	6.76	3.35	3.27
PCL [41]	1.27	5.82	2.31	2.84	3.06	3.11	1.16	2.60	7.93	3.47	3.62
PCL + prior	2.06	6.49	2.54	3.69	5.14	2.96	1.37	4.06	8.13	4.87	4.19
Ours-vgg-16 (w/o prior)	1.23	5.15	1.19	3.47	3.82	2.24	0.73	3.65	6.22	3.14	3.16
Ours-vgg-16	2.47	8.64	3.59	5.74	7.36	2.85	0.87	7.29	8.47	6.63	5.39

Table 4: mAP (%) comparison with other weakly supervised methods on EPIC KITCHENS

Methods	mAP	CorLoc
R*CNN [17]	2.54	32.68
PCL [41]	4.68	40.64
PCL + prior	6.82	46.69
Ours-vgg-16	9.75	52.53

4.4. Effect of supervision in training

Weakly supervised object detection aims to train object detection models without any bounding box labels. However, in practice it is easy and efficient to annotate at least a few bounding boxes in training images/videos. This is similar to low-shot and semi-supervised settings. We believe that it is important to test weakly supervised approaches in such a practical setting as well.

To this end, we explore the effect of adding varying amounts of ground truth object bounding box annotations into our training data. We achieve this by augmenting the losses described in Sec. 3, with an additional supervised object detection loss for videos/images where bounding box annotations are available. This loss is the same as the traditional object detection loss used in Fast-RCNN.

In practice, the IoU between object proposals and ground truth object bounding boxes is calculated and proposals having higher IoU than the threshold are considered positive samples and rest as negative. The threshold IoU is set as 0.45 to guarantee a reasonable positive samples per image. The negative and positive sample ratio is set as 5.

We compare with two baselines: (1) model without weak supervision: model trained only with supervised detection loss on images/videos with bounding box annotations and without any weakly supervised data (Ours (w/ only strong supervision)), and (2) R*CNN [17] with additional object bounding box supervision as above (R*CNN (w/ strong+weak supervision)).

We evaluate this setup on both Charades and HICO-DET datasets. The quantitative results are presented in Fig. 6. The x-axis (log scale) represents the percentage of training samples with object bounding box annotations. For example, the point $x\%$ represents that for a random $x\%$ of training data samples, the bounding box annotations are present. The remaining training samples, only have action class label. Note that 0% is the weakly-supervised setting considered earlier, while 100% represents fully-supervised setting. We observe that the mAP increases log-linearly as more su-

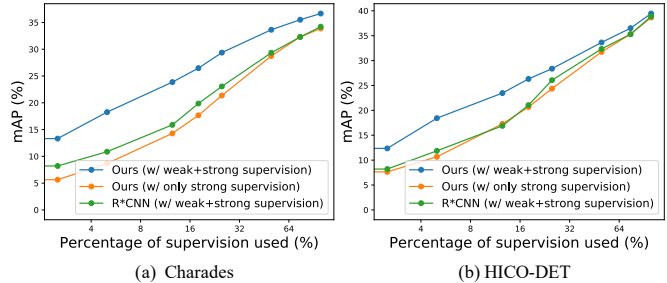


Figure 6: Performance comparison between our method trained with different supervision settings and R*CNN trained with both strong and weak supervision on (a) Charades and (b) HICO-DET.

pervision is added to the training.

For Charades, when small amount of supervision is added, we observe that our model which uses additional weakly-supervised data outperforms the model without any weak supervision. This clearly shows the potential of our weakly-supervised approach to provide complementary value in a low-shot detection setting. With as low as 70% supervision, our approach already matches the performance of fully-supervised method at 100% supervision. This means that we could cut down the amount of supervision needed to train the model without sacrificing performance. As expected, the gap between the two approaches decrease with increase in supervision. Even with “100%” bounding box annotations, our model still outperforms the fully supervised method by 2 mAP points due to joint training with action and object classification losses.

We also observe that performance gap is smaller for images (HICO-DET). We believe weak supervision is more effective in videos compared to images, where temporal linking of proposals helps in avoiding spurious detections during training.

5. Conclusion

We observe that object spatial location, appearance and movement are tightly related to the action performed with the object in images and videos. We propose a model that leverages these observations to train object detection models from samples annotated only with action labels. Comprehensive experiments are conducted on both video and image datasets. The comparison with SoTA methods shows that our approach outperforms existing weakly supervised approaches. Further our approach can also help reduce the amount of supervision required for object detection models.

References

- [1] Jean-Baptiste Alayrac, Josef Sivic, Ivan Laptev, and Simon Lacoste-Julien. Joint discovery of object states and manipulation actions. *arXiv preprint arXiv:1702.02738*, 2, 2017. 2
- [2] Peng Tang Xinggang Wang Xiang Bai and Wenyu Liu. Multiple instance detection network with online instance classifier refinement. 2017. 2, 6
- [3] Hakan Bilen and Andrea Vedaldi. Weakly supervised deep detection networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2846–2854, 2016. 2, 3, 4, 6, 7, 8
- [4] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 4724–4733. IEEE, 2017. 1
- [5] 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, 3, 5
- [6] Ramazan Gokberk Cinbis, Jakob Verbeek, and Cordelia Schmid. Weakly supervised object localization with multi-fold multiple instance learning. *IEEE transactions on pattern analysis and machine intelligence*, 39(1):189–203, 2017. 2, 3
- [7] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 5
- [8] Vincent Delaitre, Ivan Laptev, and Josef Sivic. Recognizing human actions in still images: a study of bag-of-features and part-based representations. In *BMVC 2010-21st British Machine Vision Conference*, 2010. 3
- [9] Thomas Deselaers, Bogdan Alexe, and Vittorio Ferrari. Weakly supervised localization and learning with generic knowledge. *International journal of computer vision*, 100(3):275–293, 2012. 5
- [10] Ali Diba, Vivek Sharma, Ali Mohammad Pazandeh, Hamed Pirsiavash, and Luc Van Gool. Weakly supervised cascaded convolutional networks. In *CVPR*, volume 3, page 9, 2017. 2, 6
- [11] Chuang Gan, Chen Sun, Lixin Duan, and Boqing Gong. Webly-supervised video recognition by mutually voting for relevant web images and web video frames. In *ECCV*, pages 849–866, 2016. 2
- [12] Jiyang Gao, Zhenheng Yang, Kan Chen, Chen Sun, and Ram Nevatia. Turn tap: Temporal unit regression network for temporal action proposals. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3628–3636, 2017. 2
- [13] Jiyang Gao, Zhenheng Yang, and Ram Nevatia. Cascaded boundary regression for temporal action detection. *BMVC*, 2017. 2
- [14] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 2
- [15] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition*, 2014. 2
- [16] Georgia Gkioxari, Ross Girshick, Piotr Dollár, and Kaiming He. Detecting and recognizing human-object interactions. *arXiv preprint arXiv:1704.07333*, 2017. 3
- [17] Georgia Gkioxari, Ross Girshick, and Jitendra Malik. Contextual action recognition with r* cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1080–1088, 2015. 2, 4, 6, 7, 8
- [18] Georgia Gkioxari and Jitendra Malik. Finding action tubes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 759–768, 2015. 4
- [19] Chunhui Gu, Chen Sun, Sudheendra Vijayanarasimhan, Caroline Pantofaru, David A Ross, George Toderici, Yeqing Li, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. *arXiv preprint arXiv:1705.08421*, 3(4):6, 2017. 1
- [20] Abhinav Gupta, Aniruddha Kembhavi, and Larry S Davis. Observing human-object interactions: Using spatial and functional compatibility for recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(10):1775–1789, 2009. 3
- [21] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017. 2, 3
- [22] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 2
- [23] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5
- [24] Zequn Jie, Yunchao Wei, Xiaojie Jin, Jiashi Feng, and Wei Liu. Deep self-taught learning for weakly supervised object localization. In *IEEE CVPR*, volume 2, 2017. 2
- [25] Vicky Kalogeiton, Philippe Weinzaepfel, Vittorio Ferrari, and Cordelia Schmid. Joint learning of object and action detectors. In *ICCV 2017-IEEE International Conference on Computer Vision*, 2017. 3
- [26] Vadim Kantorov, Maxime Oquab, Minsu Cho, and Ivan Laptev. Contextlocnet: Context-aware deep network models for weakly supervised localization. In *European Conference on Computer Vision*, pages 350–365. Springer, 2016. 2, 4, 6, 7
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [28] Suha Kwak, Minsu Cho, Ivan Laptev, Jean Ponce, and Cordelia Schmid. Unsupervised object discovery and tracking in video collections. In *Proceedings of the IEEE international conference on computer vision*, pages 3173–3181, 2015. 2

- [29] Subhransu Maji, Lubomir Bourdev, and Jitendra Malik. Action recognition from a distributed representation of pose and appearance. 2011. [2](#)
- [30] Seong Joon Oh, Rodrigo Benenson, Anna Khoreva, Zeynep Akata, Mario Fritz, Bernt Schiele, et al. Exploiting saliency for object segmentation from image level labels. In *IEEE Conf. Computer Vision and Pattern Recognition*, 2017. [2](#)
- [31] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. [5](#)
- [32] Julia Peyre, Ivan Laptev, Cordelia Schmid, and Josef Sivic. Weakly-supervised learning of visual relations. In *ICCV 2017-International Conference on Computer Vision 2017*, 2017. [2](#)
- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. [2](#)
- [34] 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. [3](#)
- [35] Yunhang Shen, Rongrong Ji, Changhu Wang, Xi Li, and Xuelong Li. Weakly supervised object detection via object-specific pixel gradient. *IEEE Transactions on Neural Networks and Learning Systems*, (99):1–11, 2018. [2](#)
- [36] Miaoqing Shi, Holger Caesar, and Vittorio Ferrari. Weakly supervised object localization using things and stuff transfer. In *Proc. IEEE Int. Conf. Comput. Vis.(ICCV)*, 2017. [2](#)
- [37] Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *European Conference on Computer Vision*, pages 510–526. Springer, 2016. [1](#), [2](#), [5](#)
- [38] Krishna Kumar Singh and Yong Jae Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *The IEEE International Conference on Computer Vision (ICCV)*, 2017. [2](#)
- [39] Hyun Oh Song, Ross Girshick, Stefanie Jegelka, Julien Mairal, Zaid Harchaoui, and Trevor Darrell. On learning to localize objects with minimal supervision. *arXiv preprint arXiv:1403.1024*, 2014. [2](#)
- [40] Abhilash Srikantha and Juergen Gall. Weak supervision for detecting object classes from activities. *Computer Vision and Image Understanding*, 156:138–150, 2017. [2](#)
- [41] Peng Tang, Xinggang Wang, Song Bai, Wei Shen, Xiang Bai, Wenyu Liu, and Alan Loddon Yuille. Pcl: Proposal cluster learning for weakly supervised object detection. *IEEE transactions on pattern analysis and machine intelligence*, 2018. [2](#), [6](#), [7](#), [8](#)
- [42] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013. [5](#)
- [43] Fang Wan, Pengxu Wei, Jianbin Jiao, Zhenjun Han, and Qixiang Ye. Min-entropy latent model for weakly supervised object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1297–1306, 2018. [2](#)
- [44] Limin Wang, Yuanjun Xiong, Dahua Lin, and Luc Van Gool. Untrimmednets for weakly supervised action recognition and detection. In *IEEE Conf. on Computer Vision and Pattern Recognition*, volume 2, 2017. [2](#)
- [45] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2015. [2](#)
- [46] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4724–4732, 2016. [3](#)
- [47] Yunchao Wei, Zhiqiang Shen, Bowen Cheng, Honghui Shi, Jinjun Xiong, Jiashi Feng, and Thomas Huang. Ts2c: tight box mining with surrounding segmentation context for weakly supervised object detection. In *European Conference on Computer Vision*, pages 454–470. Springer, Cham, 2018. [2](#), [6](#)
- [48] Zhenheng Yang, Jiyang Gao, and Ram Nevatia. Spatio-temporal action detection with cascade proposal and location anticipation. *BMVC*, 2017. [4](#)
- [49] Bangpeng Yao and Li Fei-Fei. Modeling mutual context of object and human pose in human-object interaction activities. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 17–24. IEEE, 2010. [3](#)
- [50] Yuan Yuan, Xiaodan Liang, Xiaolong Wang, Dit Yan Yeung, and Abhinav Gupta. Temporal dynamic graph lstm for action-driven video object detection. *ICCV*, 2017. [1](#), [2](#), [5](#), [6](#), [7](#)
- [51] Xiaopeng Zhang, Jiashi Feng, Hongkai Xiong, and Qi Tian. Zigzag learning for weakly supervised object detection. *arXiv preprint arXiv:1804.09466*, 2018. [2](#)
- [52] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas Huang. Adversarial complementary learning for weakly supervised object localization. In *IEEE CVPR*, 2018. [2](#)
- [53] Yongqiang Zhang, Yancheng Bai, Mingli Ding, Yongqiang Li, and Bernard Ghanem. W2f: A weakly-supervised to fully-supervised framework for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 928–936, 2018. [2](#)
- [54] Luowei 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. [2](#)
- [55] Yi Zhu, Yanzhao Zhou, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. Soft proposal networks for weakly supervised object localization. In *Proc. IEEE Int. Conf. Comput. Vis.(ICCV)*, pages 1841–1850, 2017. [2](#)