

Image Co-localization by Mimicking a Good Detector's Confidence Score Distribution

Yao Li, Lingqiao Liu, Chunhua Shen^(✉), and Anton van den Hengel

School of Computer Science, The University of Adelaide, Adelaide, Australia
`chunhua.shen@adelaide.edu.au`

Abstract. Given a set of images containing objects from the same category, the task of image co-localization is to identify and localize each instance. This paper shows that this problem can be solved by a simple but intriguing idea, that is, a common object detector can be learnt by making its detection confidence scores distributed like those of a strongly supervised detector. More specifically, we observe that given a set of object proposals extracted from an image that contains the object of interest, an accurate strongly supervised object detector should give high scores to only a small minority of proposals, and low scores to most of them. Thus, we devise an entropy-based objective function to enforce the above property when learning the common object detector. Once the detector is learnt, we resort to a segmentation approach to refine the localization. We show that despite its simplicity, our approach outperforms state-of-the-arts.

Keywords: Image co-localization · Unsupervised object discovery

1 Introduction

There has been an explosion of images available on the Internet in recent years, largely due to the popularity of photo sharing sites like Facebook and Flickr. However, most of these images are either unlabeled or weakly-labeled. One way of accessing these images is finding images depicting the same object, for instance, Google Image Search will return images containing a common object described by the user input keyword. In this paper, we aim to localize the common object in this scenario (without using any other forms of supervision, *e.g.*, manually-labeled negative samples). This task is known as the image co-localization task in literature [4, 17, 30].

Image co-localization is a particularly challenging task, and thus there exist a limited number of comparable methods [4, 17, 30]. These methods address this problem from various perspectives. The work in [30] introduces binary latent variables to indicate the presence of the common object and formulates the co-localization via latent variable inference. The work of [4], in contrast, localizes the common object by matching common object parts. Our work differs from

First two authors contributed equally.

previous approaches in that it directly learns the common object detector by modeling its detection confidence score distribution on each image, and achieves the localization with the learned detector.

The key insight of our method is that although we do not have sufficient supervision to learn a strongly supervised object detector, it is still possible to learn an “artificial” detector by modeling its detection confidence score distribution on object proposals [31, 34]. The intuition is inspired from the behaviour of an accurate strongly supervised object detector, that is, when applied to object proposals extracted from an image contains the object of interest, only a small minority of proposals will be given high detection confidence scores while most of them are associated with low scores. Motivated the above observation, in this paper we design a novel Shannon-entropy-based objective function to promote the scarcity of high detection confidence scores within an image while avoiding the trivial solution of producing low scores for all proposals. In other words, by optimizing the proposed objective, our approach will encourage the existence of a few high response proposals in each image as the common object while suppressing responses in the remainder proposals which will be deemed as background.

To generate the final co-localization results, we have also devised a method for improving the bounding box estimate. Inspired by detection-by-segmentation approaches (*e.g.*, [22]), we use the final detection heat map and color information to define a CRF-based segmentation algorithm, the output of which indicates the instances of the common object.

Through an extensive evaluation on several benchmark datasets, including the PASCAL VOC 2007 and 2012 [8], and also some subsets of the ImageNet [6], we demonstrate that our approach outperforms the state-of-the-arts for the image co-localization task.

2 Related Work

Image co-localization shares some similarities with image co-segmentation [3, 16, 26] in the sense that both problems require a set of images containing objects from a common category as input. Instead of generating a precise segmentation of the related objects in each image, co-localization algorithms [4, 17, 30] aim to draw a tight bounding box around the object. Image co-localization is also related to works on weakly supervised object localization (WSOL) [1, 5, 7, 24, 28, 29, 32, 33] as both try to localize objects of the same type within an image set, the key difference is WSOL requires manually-labeled negative images whereas co-localization does not.

Tang *et al.* [30] formulate co-localization as a boolean constrained quadratic program which can be relaxed to a convex problem, which is further accelerated by the Frank-Wolfe Algorithm [17]. Recently, Cho *et al.* [4] propose a Probabilistic Hough Matching algorithm to match object proposals across images and then dominant objects are localized by selecting proposals based on matching scores. There are also approaches address the problem of co-localization in video [17, 21, 23]. Notably, Prest *et al.* [23] select spatio-temporal tubes which

are likely to contain the common object, and Joulin *et al.* [17], in contrast, extend [30] by incorporating temporal consistency.

However, in this paper, we tackle the co-localization problem from a new perspective, that is, learning the common object detector by modeling its detection confidence score distribution, and thus get rid of the need of manually-labeled negative images. An advantage of the proposed approach for learning common object detectors is that it provides an explicit mechanism by which to exploit the relationship between localization and detection. The benefits of exploiting this relationship have been identified before in WSOL. For example, in [7], objects are localized by minimizing a Conditional Random Field (CRF) energy function which incorporates class-specific information, and the class-specific information is learned from the localized objects. Cinbis *et al.* [5] propose a multi-fold training procedure for Multiple Instance Learning whereby, at each iteration, positive instances in each fold are localized by a detector trained from other folds in the previous iteration. The approach that we propose here, however, is the first to systematically leverage the idea of jointly performing object detection and localization for co-localizing common objects in images.

3 Approach

We give an overview of our image co-localization framework in Fig. 1. The input to our framework is a set of N images $\mathcal{I} = \{I_1, I_2, \dots, I_N\}$ contains one common object (*e.g.*, aeroplane), and we aim to annotate the location of common object

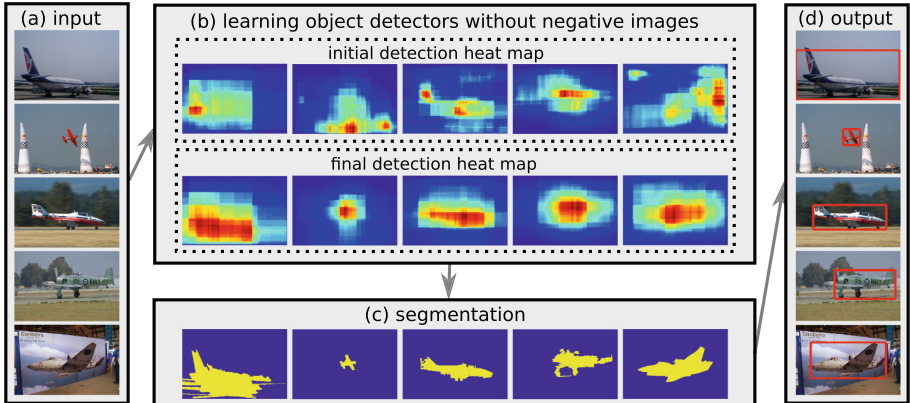


Fig. 1. An overview of our image co-localization framework. (a) The input of our system is a set of images contains a common object category (here, aeroplane). (b) The common object detector is learnt by modeling the distribution of detection confidence scores. (c) Detection heat maps generated by the learnt detector are used as the unary potential for graph-cuts segmentation. (d) The output for each image is the smallest rectangle which covers the corresponding segmentation.

instances in each image. Inspired by the behaviour of an accurate strongly supervised object detector (Sect. 3.1), the core of our framework is the procedure of learning the common object detector by modeling its detection confidence score distribution (Sect. 3.2). We further formulate object localization as a segmentation problem (Sect. 3.3), which involves using the detection heat map to define unary potentials of a binary energy function and solving it efficiently by standard graph-cuts.

3.1 The Behaviour of an Accurate Strongly Supervised Detector

Object proposals [31, 34], which are image regions that are likely to contain objects, have been widely used in recent object detection approaches [10–12]. In this section we are interested in the statistics of proposal detection confidence scores on an image generated by a strongly supervised detector. The observation here motivates our formulation for learning common object detectors in Sect. 3.2.

More specifically, we apply one state-of-the-art strongly supervised object detector Fast R-CNN [10] (trained on PASCAL VOC 2007 trainval set [8]) to a PASCAL VOC 2007 test image which contains the object of interest (Fig. 2(a)). After obtaining the detection confidence scores of the more than 2000 object proposals [31] extracted from this image, we calculate the normalized histogram of detection confidence scores of all proposals (Fig. 2(b)).

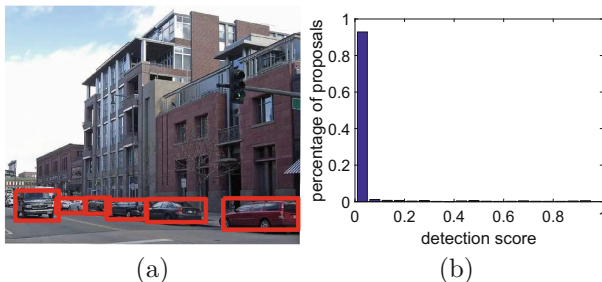


Fig. 2. (a) Predicted objects by Fast R-CNN [10]. (b) Normalized detection confidence score histogram of object proposals in (a). We observe the same statistics for most images.

From Fig. 2(b) it is clear that, although there are multiple instances of the object of interest (“car” in this case), more than 90 % of object proposals have a very low detection confidence score (less than 0.05), which indicates that a dominantly large portion of proposals are likely to cover image regions that do not cover the object of interest tightly. This is understandable as object proposal generation is a pre-processing step in object detection systems, where recall rate is much more important than precision (not missing any objects of interest is more important than generating less false positives).

3.2 Learning Detectors by Modeling Detection Score Distribution

In the setting of image co-localization, although all we know is that there exists a common object category across images, we still aim to learn the common object detector. This is possible by modeling the distribution of proposals detection confidence scores. More specifically, in our method the common object detector will be learned by enforcing its the distribution of detection confidence scores to mimic that of an accurate strongly supervised detector (Sect. 3.1).

Formally, for each image $I_i \in \mathcal{I}$, we first extract a set of object proposals $\mathcal{B}_i = \{B_{i,1}, B_{i,2}, \dots, B_{i,M_i}\}$ using EdgeBox [34], which shows good performance in a recent review [14]. Let $\phi(B_{i,j}) \in \mathbb{R}^K$ denote the feature representation of proposal $B_{i,j} \in \mathcal{B}_i$. The particular detection confidence scores that we use are formulated as follows

$$s_{i,j} = f(\mathbf{w}^T \phi(B_{i,j}) + b), \quad (1)$$

where $\mathbf{w} \in \mathbb{R}^K$, $b \in \mathbb{R}^1$ denote weight and bias terms of the detector respectively, and $f(\cdot)$ is the softplus function which has the form $f(x) = \ln(1 + \exp(x))$.

Irrespective of the form of the detector, we can construct the set of detection confidence scores $\mathcal{S}_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,M_i}\}$ over all the proposals \mathcal{B}_i , and normalize them as $p_{i,j} = \frac{s_{i,j} + \epsilon}{\sum_j (s_{i,j} + \epsilon)}$, where the parameter ϵ is a small constant. If the detector in Eq. (1) is trained with strong supervision, according to the observation in Sect. 3.1, most of its detection confidence scores in \mathcal{S}_i should have near-zero values which means that the score vector $\mathbf{s}_i = [s_{i,1}, s_{i,2}, \dots, s_{i,M_i}]^T$ and its normalized version $\mathbf{p}_i = [p_{i,1}, p_{i,2}, \dots, p_{i,M_i}]^T$ should be sparse vectors. Note that when all proposals have zero detection confidence scores, \mathbf{s}_i will be sparse but \mathbf{p}_i will be dense due to the effect of the constant ϵ . Thus, our method will be based on \mathbf{p}_i because enforcing its sparsity will be equivalent to requiring the detector to have few high detection confidence scores and many low (zero) detection confidence scores, in other words, the detection confidence score distribution will mimic that of an accurate strongly supervised detector.

Objective Function. To measure the sparsity of the normalized detection confidence score vector \mathbf{p}_i , we opt for the Shannon entropy in this work, that is,

$$\mathcal{L}(\mathbf{p}_i) = - \sum_{j=1}^{M_i} p_{i,j} \log p_{i,j}, \quad (2)$$

and the objective for learning the common object detector is formulated as follows:

$$\min_{\mathbf{w}, b} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\mathbf{p}_i) + \lambda \|\mathbf{w}\|_2^2, \quad (3)$$

where we use the square of the L_2 -norm of \mathbf{w} as a regularizer on the weight vector.

So the optimal value of the weight and bias of the detector is given by:

$$\mathbf{w}^*, b^* = \operatorname{argmin}_{\mathbf{w}, b} - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{M_i} p_{i,j} \log p_{i,j} + \lambda \|\mathbf{w}\|_2^2. \quad (4)$$

Note that Eq. (4) does not involve a set of manually-labeled negative samples which do not contain the object of interest, but rather describes the desired form of the detection confidence score distribution of object proposals. The learning process also implicitly takes advantage of the chicken-and-egg relationship between object localization and detection: precisely localized object instances are critical for training a good object detector, and objects can be localized more precisely by a well-trained detector.

Optimization. As our objective function in Eq. (4) is non-convex, we minimize it using stochastic gradient descent (SGD). Similar to the approach used in training a Convolutional Neural Network [19], we divide all data (*i.e.*, object proposals) into mini-batches. We initialize the weight vector \mathbf{w} from a zero-mean Gaussian distribution, while the bias term b is set to zero initially. During training we divide the learning rate (which is set to 0.1 initially) by 10 after each 10 epochs. We stop learning after 20 epochs when the objective function converges.

Modification. After minimizing Eq. (4), when we visualize the proposal with the maximal detection confidence score for each image (Fig. 3), it is interesting to note that the learnt detector may not fire at the common object but some common visual patterns (*e.g.*, common object parts, common object with some context) instead. Also, the discovered common visual patterns can be very different if the initialization of our detector varies (different local minimums). However in this work, as we aim to co-localize the common object, we reformulate Eq. (1) by incorporating the “objectness” score $o_{i,j}$ (outputs of Edgebox) of each proposal $B_{i,j}$ as a weight to favour proposals with high objectness score (which more likely to cover a whole object tightly)

$$s_{i,j} = o_{i,j} f(\mathbf{w}^T \phi(B_{i,j}) + b). \quad (5)$$

We experimentally find that minimizing Eq. (4) using $s_{i,j}$ defined in Eq. (5) gives a stable solution regardless of initialization.

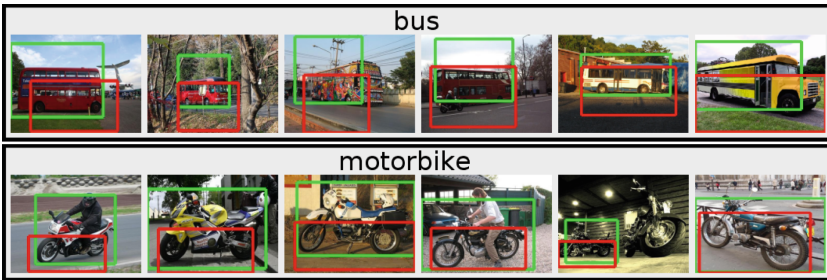


Fig. 3. Our detectors fire at different common visual patterns (denoted by red and green bounding boxes) by minimizing Eq. (4) with different random initializations. Although these common visual patterns may not be suitable for the co-localization task, they may be useful for other computer vision tasks, such as discovering common object parts for fine-grained image classification [18]. (Color figure online)

Localizing the Common Object. The optimal \mathbf{w} and b , inserted into Eq. (1), lead to a mechanism for determining the detection confidence scores for all object proposals. The nature of the co-localization problem means that the maximal score for each image indicates the desired detection. This method is used as a baseline in the Experiments section (Sect. 4.1).

Discussion. Theoretically, other sparsity measures could be employed to replace the Shannon entropy. Note that the commonly used L_1 norm cannot be applied here because $\|\mathbf{p}_i\|_1 = 1$. One possible way to use L_1 norm is to redefine the normalization score $p_{i,j} = \frac{s_{i,j} + \epsilon}{\sqrt{\sum_j (s_{i,j} + \epsilon)^2}}$.

3.3 Refining the Bounding Box Estimate

The quality of the detections generated through the above described process depends entirely on the quality of object proposals. To overcome this dependency, and enable better final bounding box estimates to be achieved, we have developed a bounding box refinement process as follows.

Given the optimal \mathbf{w}^* and b^* identified by minimizing Eq. (4), we generate the detection heat map as follows. For each pixel in the image, we add up the weighted detection confidence score $s_{i,j}$ from Eq. (5) for all proposals $B_{i,j}$ that cover this pixel (zero for pixels not covered by any proposals). The values are then normalized to the interval $[0, 1]$. This gives rise to a set of detection heat maps $\mathcal{H} = \{H_1, H_2, \dots, H_n\}$. Some examples are illustrated in Fig. 4.

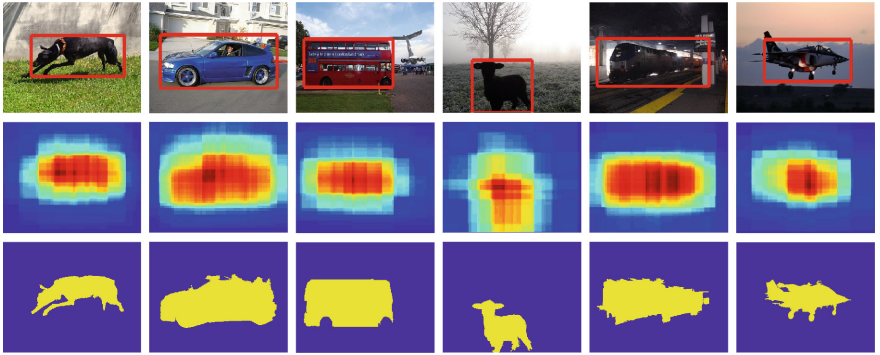


Fig. 4. Examples of our co-localization process. From top to bottom: input images (predicted boxes in red), detection heat maps, segmentation results. (Color figure online)

Given the set of detection heat maps \mathcal{H} , we aim to produce a segmentation of the entire object. This approach is inspired by previous work which casts localization as a segmentation problem (*e.g.*, [22]).

Formally, we formulate the segmentation problem as a standard graph-cut problem. We first extract superpixels [9] to construct the vertex set $\{m\}$ and

aim to label each superpixel as foreground ($y_m = 1$) or background ($y_m = 0$). Mathematically, the energy function is given by

$$E(\mathbf{y}) = \sum_m u_m(y_m) + \sum_{(m,n) \in \mathcal{E}} v_{mn}(y_m, y_n), \quad (6)$$

where u_m and v_{mn} are the unary and pairwise potential respectively. \mathcal{E} is the set of edges connecting superpixels¹.

Unary Potential u_m . Inspired by [20], the unary potential is the novel part of our segmentation framework, which carries information from the detection heat map H :

$$u_p(y_m) = -\log A_m(y_m), \quad (7)$$

where A_m is the prior information from the detection heat map H :

$$\begin{aligned} A_p(y_m = 1) &= H(m), \\ A_p(y_m = 0) &= 1 - H(m), \end{aligned} \quad (8)$$

where $H(m)$ is the mean of values inside superpixel m on map H .

Pairwise Potential v_{mn} . Our pairwise potential is defined as follows.

$$v_{mn}(y_m, y_n) = [y_m \neq y_n] e^{-\beta \|C(m) - C(n)\|_2^2}, \quad (9)$$

where $C(m)$ is the color histogram feature. As in [20, 25], this potential penalizes superpixels with different colors taking the same label.

As our pairwise potential in Eq. (9) is submodular, the optimal label \mathbf{y}^* can be found efficiently by the graph-cuts [2]. As shown in Fig. 4, the segmentation derived through this approach is accurate. The final bounding box estimate is then calculated as the smallest rectangle which covers the segmentation.

4 Experiments

Datasets. We evaluate our approach on three datasets, including VOC 2007 and 2012 [8] datasets, six subsets of the ImageNet dataset [6] which have not been used in the ILSVRC [27]². For VOC datasets, following previous works in co-localization and weakly supervised object localization [1, 4, 5, 33], we use all images on the *trainval* set discarding images that only contain object instances marked as “difficult” or “truncate”.

Evaluation Metric. We use two metrics to evaluate our approach. Firstly, for comparison with state-of-the-art approaches, we use the CorLoc metric [7],

¹ In our case two superpixels are connected if the distance between their centroids is smaller than the sum of their major axis length.

² The six categories are chipmunk, rhino, stoat, racoon, rake and wheelchair. Note that ground-truth bounding box annotations are available for these categories, thus enable quantitative evaluation.

which is defined as the percentage of images that are correctly localized. An image is considered as correctly localized if the Intersection-over-Union (IoU) score between the predicted bounding box and any ground-truth bounding boxes of the object of interest exceeds 50 %.

Implementation Details. We use Edgebox [34] to extract object proposals with a maximum of 2000 proposals extracted from each image. We represent each Edgebox proposal as a 4096-dimensional CNN feature from the *fc6* layer (after ReLU) from the *BVLC Reference CaffeNet* model [15]. We use a fixed value of 1 for λ in Eq. (4) which controls the tradeoff between the loss function and regularizer. The value of β in Eq. (9) is set to 10.

4.1 Ablation Study

Baselines. To investigate the impact of the various elements of our approach, we consider the following two baseline methods:

- “obj-sel”: the predicted bounding box for an image is simply the proposal with maximum objectness score.
- “obj-seg”: for each image, objectness scores of all proposals are treated as detection confidence scores to generate a fake detection heat map, which is then sent to our segmentation model in Sect. 3.3.

The two methods proposed in our work are:

- “our-sel”: given the learnt detector (Sect. 3.2), we simply select the object proposal which has the maximum detection confidence score $p_{i,j}$ i.e., $B_i^* = \operatorname{argmax}_{B_i, j \in \mathcal{B}_i} S_{i,j}$.
- “our-seg”: combination of detector training (Sect. 3.2) and segmentation refinement (Sect. 3.3).

Corloc scores for the above four methods on the VOC 2007 dataset are illustrated in Fig. 5.

As shown in Fig. 5, the simplest baseline “obj-sel” does not work well (19.8 % CorLoc). This is because the objectness measure of Edgebox [34] is heuristically defined based on only edge information, which does not exploit the common object assumption.

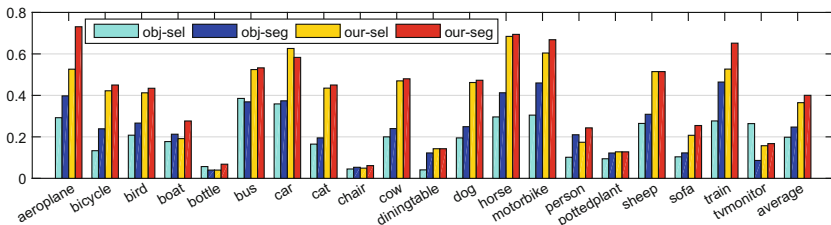


Fig. 5. CorLoc scores of our approaches, and baselines, on the VOC 2007 dataset.

However, the “obj-seg” baseline in which we use objectness scores to generate a detection heat map for each image, performs quite well, with CorLoc increasing to 24.7%. Surprisingly, this performance is on the par with one state-of-the-art image co-localization approach [17] (24.7% *vs.* 24.6%), even though there is no common object assumption. This phenomenon indicates that our segmentation model is quite effective.

Thanks to our common object detector learning procedure in Sect. 3.2, “our-sel” achieves a performance of 36.5%, outperforming “obj-sel” and “obj-seg” by over 16% and 11% respectively. This verifies the effectiveness of this procedure, and particularly that, although we do not have annotated image labels nor bounding boxes, the detector still captures the appearance of the common object, which improves co-localization significantly.

Combing the advantages of the common object detector learning procedure (Sect. 3.2) and segmentation refinement (Sect. 3.3), we observe another 3.5% boost in the case of “our-seg”, reaching 40.0% Corloc. Thus we use “our-seg” to compare with state-of-the-art approaches.

Number of Candidate Proposals. To evaluate the robustness of our approach under different number of candidate object proposals, we test three settings—500, 1000 and 2000, which results in 39.2%, 39.6% and 40.0% CorLoc respectively. This indicates that our approach is quite insensitive to the changes in the number of candidate proposals.

4.2 Diagnosing the Localization Error

In order to better understand the localization errors, following [5, 13], each predicted bounding box predicted by our approach is categorized into the following five cases: (1) correct: IoU score exceeds 50%, (2) g.t. in hypothesis: ground-truth completely inside prediction, (3) hypothesis in g.t.: prediction completely inside ground-truth, (4) no overlap: IoU score equals zero, (5) low overlap: none of the above four cases. In Fig. 6 we show the error modes of our approach across all categories on the VOC 2007 dataset.

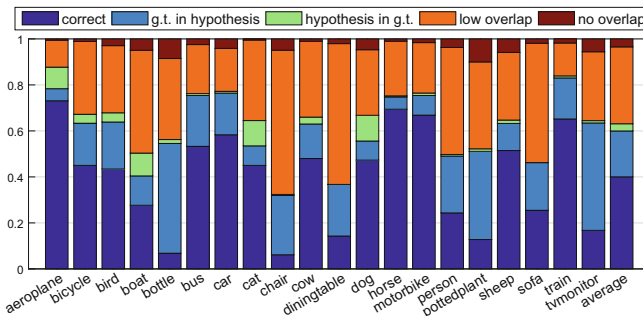


Fig. 6. An illustration of error types for our approach on the VOC 2007 dataset.

As shown in Fig. 6, the fraction of “no overlap” cases is quite small (3.5%) across all categories, which means our approach can localize common objects to some extent in most cases. Comparing “g.t. in hypothesis” to its “hypothesis in g.t.”, it is clear that the former appears more frequently (19.9% *v.s.* 3.1%), which means our approach tends to localize objects with some context details. In terms of correct localization, the three categories with lowest CorLoc values are *bottle* (6.8%), *chair* (6.2%) and *pottedplant* (12.8%). Objects in these categories are always in very clustered environments with occlusion (*e.g.*, chair is often occluded by table) which makes the task quite challenging.

4.3 Comparison to State-of-the-Art Approaches

Comparison to Image Co-localization Approaches. We now compare the results of our approach to the state-of-the-art image co-localization approaches of Joulin *et al.* [17] and Cho *et al.* [4] on the VOC 2007 dataset (Table 1). The performance of our approach exceeds that of Joulin *et al.* [17] significantly in most categories, with an improvement of over 15% in mean CorLoc. The recent approach of Cho *et al.* [4] relies on matching object parts by Hough Transform with the predicted bounding box is selected by a heuristic standout score. Candidate regions are object proposals represented by whitened HOG features. However, we found that this whitening process, whose mean vector and covariance matrix are estimated from the random sampled images from the same dataset (inevitably using images from other categories), is crucial for the performance of their algorithm. Our performance bypasses that of [4] by a reasonable margin of 3.4%.

Table 1. Comparison to image co-localization approaches on the VOC 2007 dataset in terms of CorLoc metric [7].

VOC	Aero	Bike	Bird	Boat	Bottle	Bus	Car	Cat	Chair	Cow	Table	Dog	Horse	Mbike	Person	Plant	Sheep	Sofa	Train	Tv	Mean
[17]	32.8	17.3	20.9	18.2	4.5	26.9	32.7	41.0	5.8	29.1	34.5	31.6	26.1	40.4	17.9	11.8	25.0	27.5	35.6	12.1	24.6
[4]	50.3	42.8	30.0	18.5	4.0	62.3	64.5	42.5	8.6	49.0	12.2	44.0	64.1	57.2	15.3	9.4	30.9	34.0	61.6	31.5	36.6
Ours	73.1	45.0	43.4	27.7	6.8	53.3	58.3	45.0	6.2	48.0	14.3	47.3	69.4	66.8	24.3	12.8	51.5	25.5	65.2	16.8	40.0

To further verify the effectiveness of our approach, we now present an evaluation on the VOC 2012 dataset [8] which has twice the number of images of VOC 2007. Table 2 shows our performance along with that of Cho *et al.* [4] which we evaluated using their publicly available code. It is clear that on average our approach outperforms that of Cho *et al.* [4] by 2%.

Table 2. Comparison to image co-localization approaches on the VOC 2012 dataset in terms of CorLoc metric [7].

VOC	Aero	Bike	Bird	Boat	Bottle	Bus	Car	Cat	Chair	Cow	Table	Dog	Horse	Mbike	Person	Plant	Sheep	Sofa	Train	Tv	Mean
[4]	57.0	41.2	36.0	26.9	5.0	81.1	54.6	50.9	18.2	54.0	31.2	44.9	61.8	48.0	13.0	11.7	51.4	45.3	64.6	39.2	41.8
Ours	65.7	57.8	47.9	28.9	6.0	74.9	48.4	48.4	14.6	54.4	23.9	50.2	69.9	68.4	24.0	14.2	52.7	30.9	72.4	21.6	43.8

Comparison to Weakly Supervised Object Localization Approaches.

We also compare our approach with some state-of-the-art approaches on weakly supervised object localization. Table 3 illustrates the comparison of several recent works and our approach on VOC 2007 dataset. In particular, our performance (40.0%) is comparable to that of a very recent work [33] (40.2%) which also uses CNN features and Edgebox proposals. As shown in Table 3, though we do not have any negative images, we still outperforms WSOL approaches on 3 of 20 categories.

Table 3. Comparison to weakly supervised object localization approaches on the VOC 2007 dataset in terms of CorLoc metric [7]. Note that these comparators require access to a negative image set, whereas our approach does not.

VOC	Aero	Bike	Bird	Boat	Bottle	Bus	Car	Cat	Chair	Cow	Table	Dog	Horse	Mbike	Person	Plant	Sheep	Sofa	Train	Tv	Mean
[29]	42.4	46.5	18.2	8.8	2.9	40.9	73.2	44.8	5.4	30.5	19.0	34.0	48.8	65.3	8.2	9.4	16.7	32.3	54.8	5.5	30.4
[28]	67.3	54.4	34.3	17.8	1.3	46.6	60.7	68.9	2.5	32.4	16.2	58.9	51.5	64.6	18.2	3.1	20.9	34.7	63.4	5.9	36.2
[5]	56.6	58.3	28.4	20.7	6.8	54.9	69.1	20.8	9.2	50.5	10.2	29.0	58.0	64.9	36.7	18.7	56.5	13.2	54.9	59.4	38.8
[33]	37.7	58.8	39.0	4.7	4.0	48.4	70.0	63.7	9.0	54.2	33.3	37.4	61.6	57.6	30.1	31.7	32.4	52.8	49.0	27.8	40.2
[1]	66.4	59.3	42.7	20.4	21.3	63.4	74.3	59.6	21.1	58.2	14.0	38.5	49.5	60.0	19.8	39.2	41.7	30.1	50.2	44.1	43.7
[24]	79.2	56.9	46.0	12.2	15.7	58.4	71.4	48.6	7.2	69.9	16.7	47.4	44.2	75.5	41.2	39.6	47.4	32.2	49.8	18.6	43.9
[32]	80.1	63.9	51.5	14.9	21.0	55.7	74.2	43.5	26.2	53.4	16.3	56.7	58.3	69.5	14.1	38.3	58.8	47.2	49.1	60.9	48.5
Ours	73.1	45.0	43.4	27.7	6.8	53.3	58.3	45.0	6.2	48.0	14.3	47.3	69.4	66.8	24.3	12.8	51.5	25.5	65.2	16.8	40.0

We have also conducted an object detection experiment on VOC 2007. Specifically, for each category, we treated predicted bounding boxes of our co-localization algorithm on trainval set as ground-truth annotations and sampled proposals from other categories or have a overlap ratio less than 0.1 against our localized bounding boxes as negative samples. The fc6 feature from the CaffeNet are extracted and hard negative mining is performed to train the detector. We achieve a mAP of 16.7% on the testset when using a nms threshold of 0.5. Although our performance is lower than some WSOL approaches, it is understandable as we do not use negative data for co-localization. Moreover, we can easily extend our formulation (Eq. 4) to handle negative data and thus perform WSOL.

Visualization. In Fig. 7, we provide a set of successful co-localization results along with the corresponding detection heat maps for some categories of the VOC 2007 dataset. It demonstrates that detection heat maps successfully predict the correct location of the common object regardless of changes in scale, appearance and viewpoint. This provides a strong indication that, although trained without annotated positive or negative examples, our approach is able to discriminate the common object from other objects in the scene.

4.4 ImageNet Subsets

We note that the CNN model used for extracting features is pre-trained in the ILSVRC [27], whose training set may have some overlapping categories with the VOC datasets. In order to justify our approach is insensitive to the object

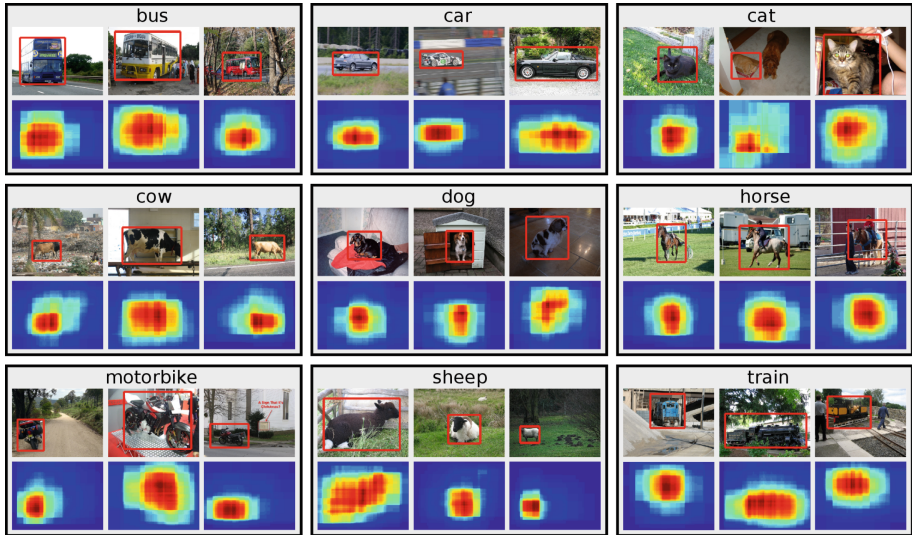


Fig. 7. Examples of successful co-localization results for the VOC 2007 dataset. For each category, the top row depicts predicted bounding boxes on the original image, the bottom row shows corresponding detection heat maps.

category, we randomly selected six subsets of the ImageNet [6] which have not been used in the ILSVRC (thus “unseen” by the CNN model) for evaluation.

Table 4 shows our co-localization result along with that of the current state-of-the-art work of Cho *et al.* [4]. Clearly, our approach outperforms [4] by a reasonable margin on all categories except the *rhino* category, whose images tend to have relatively large common instances and less cluttered background. Some successfully co-localization samples are depicted in Fig. 8.

Table 4. Comparison to image co-localization approaches on the ImageNet subsets in terms of CorLoc metric [7]. Note that these categories have not been used for pre-training the CNN model, which is used as a feature extractor in this work.

ImageNet	Chipmunk	Rhino	Stoat	Raccoon	Rake	Wheelchair	Mean
Cho <i>et al.</i> [4]	26.6	81.8	44.2	30.1	8.3	35.3	37.7
Ours	44.9	81.8	67.3	41.8	14.5	39.3	48.3

We also visualize some failure cases of the two categories our approach performed worst—*rake* and *wheelchair* (Fig. 9). Interestingly, these failure cases are quite understandable. For example, a large portion of images in the *rake* category contains both people and rakes, thus our approach tends to capture this combination as the “common object”. A similar phenomenon is also observed in the *wheelchair* category in which people occur along with wheelchairs.

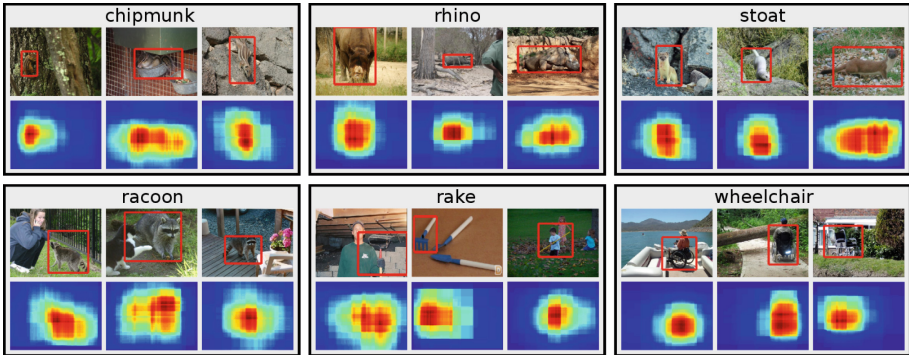


Fig. 8. Examples of successful co-localization results for the ImageNet subsets.

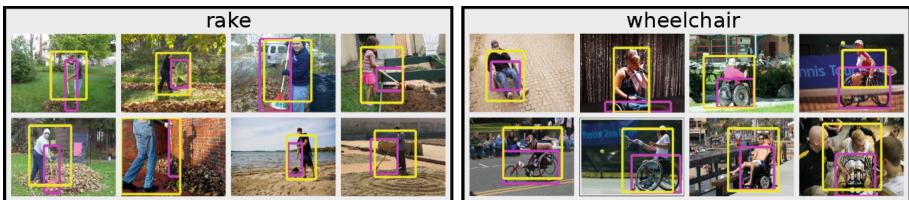


Fig. 9. Examples of failure cases of *rake* and *wheelchair* (ground truth in megenta and predicted boxes in yellow). (Color figure online)

5 Conclusion

We have addressed the image co-localization problem by directly learning a common object detector. The key discovery made in this paper is that this detector can be learned with the objective of making its detection score distribution mimic an accurate strongly supervised object detector. Also, we have illustrated that it is profitable to use a CRF model to refine the co-localization result, which has not been explored in recent works on co-localization.

Acknowledgements. This work was in part supported by the Data to Decisions CRC Centre. C. Shen’s participation was in part supported by ARC Future Fellowship No. FT120100969.

References

1. Bilen, H., Pedersoli, M., Tuytelaars, T.: Weakly supervised object detection with convex clustering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1081–1089 (2015)
2. Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach. Intell.* **23**(11), 1222–1239 (2001)

3. Chen, X., Shrivastava, A., Gupta, A.: Enriching visual knowledge bases via object discovery and segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2035–2042 (2014)
4. Cho, M., Kwak, S., Schmid, C., Ponce, J.: Unsupervised object discovery and localization in the wild: part-based matching with bottom-up region proposals. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1201–1210 (2015)
5. Cinbis, R.G., Verbeek, J.J., Schmid, C.: Multi-fold MIL training for weakly supervised object localization. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2409–2416 (2014)
6. Deng, J., Dong, W., Socher, R., Li, L., Li, K., Li, F.: ImageNet: a large-scale hierarchical image database. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255 (2009)
7. Deselaers, T., Alexe, B., Ferrari, V.: Weakly supervised localization and learning with generic knowledge. *Int. J. Comput. Vis.* **100**(3), 275–293 (2012)
8. Everingham, M., Eslami, S.M.A., Gool, L.V., Williams, C.K.I., Winn, J.M., Zisserman, A.: The PASCAL visual object classes challenge: a retrospective. *Int. J. Comput. Vis.* **111**(1), 98–136 (2015)
9. Felzenszwalb, P.F., Huttenlocher, D.P.: Efficient graph-based image segmentation. *Int. J. Comput. Vis.* **59**(2), 167–181 (2004)
10. Girshick, R.: Fast R-CNN. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1440–1448 (2015)
11. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Region-based convolutional networks for accurate object detection and segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* **38**(1), 142–158 (2016)
12. He, K., Zhang, X., Ren, S., Sun, J.: Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* **37**(9), 1904–1916 (2015)
13. Hoiem, D., Chodpathumwan, Y., Dai, Q.: Diagnosing error in object detectors. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) *ECCV 2012, Part III. LNCS*, vol. 7574, pp. 340–353. Springer, Heidelberg (2012)
14. Hosang, J.H., Benenson, R., Dollár, P., Schiele, B.: What makes for effective detection proposals? *IEEE Trans. Pattern Anal. Mach. Intell.* **38**(4), 814–830 (2016)
15. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: convolutional architecture for fast feature embedding. *arXiv preprint [arXiv:1408.5093](https://arxiv.org/abs/1408.5093)* (2014)
16. Joulin, A., Bach, F.R., Ponce, J.: Discriminative clustering for image co-segmentation. In: *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition*, pp. 1943–1950 (2010)
17. Joulin, A., Tang, K., Fei-Fei, L.: Efficient image and video co-localization with Frank-Wolfe algorithm. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014, Part VI. LNCS*, vol. 8694, pp. 253–268. Springer, Heidelberg (2014)
18. Krause, J., Jin, H., Yang, J., Li, F.: Fine-grained recognition without part annotations. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5546–5555 (2015)
19. Krizhevsky, A., Sutskever, I., Hinton, G.E.: ImageNet classification with deep convolutional neural networks. In: *Proceedings of the Advances in Neural Information Processing Systems*, pp. 1106–1114 (2012)
20. Küttel, D., Ferrari, V.: Figure-ground segmentation by transferring window masks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 558–565 (2012)

21. Kwak, S., Cho, M., Ponce, J., Schmid, C., Laptev, I.: Unsupervised object discovery and tracking in video collections. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3173–3181 (2015)
22. Parkhi, O.M., Vedaldi, A., Jawahar, C.V., Zisserman, A.: The truth about cats and dogs. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1427–1434 (2011)
23. Prest, A., Leistner, C., Civera, J., Schmid, C., Ferrari, V.: Learning object class detectors from weakly annotated video. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3282–3289 (2012)
24. Ren, W., Huang, K., Tao, D., Tan, T.: Weakly supervised large scale object localization with multiple instance learning and bag splitting. *IEEE Trans. Pattern Anal. Mach. Intell.* **38**(2), 405–416 (2016)
25. Rother, C., Kolmogorov, V., Blake, A.: GrabCut: interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.* **23**(3), 309–314 (2004)
26. Rubinstein, M., Joulin, A., Kopf, J., Liu, C.: Unsupervised joint object discovery and segmentation in internet images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1939–1946 (2013)
27. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.S., Berg, A.C., Li, F.: ImageNet large scale visual recognition challenge. *Int. J. Comput. Vis.* **115**(3), 211–252 (2015)
28. Shi, Z., Hospedales, T.M., Xiang, T.: Bayesian joint topic modelling for weakly supervised object localisation. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2984–2991 (2013)
29. Siva, P., Xiang, T.: Weakly supervised object detector learning with model drift detection. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 343–350 (2011)
30. Tang, K., Joulin, A., Li, L., Li, F.: Co-localization in real-world images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1464–1471 (2014)
31. Uijlings, J.R.R., van de Sande, K.E.A., Gevers, T., Smeulders, A.W.M.: Selective search for object recognition. *Int. J. Comput. Vis.* **104**(2), 154–171 (2013)
32. Galleguillos, C., Babenko, B., Rabinovich, A., Belongie, S., Wang, C., Ren, W., Huang, K., Tan, T.: Weakly supervised object localization with latent category learning. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014, Part VI. LNCS*, vol. 8694, pp. 431–445. Springer, Heidelberg (2014)
33. Wang, X., Zhu, Z., Yao, C., Bai, X.: Relaxed multiple-instance SVM with application to object discovery. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1224–1232 (2015)
34. Zitnick, C.L., Dollár, P.: Edge boxes: locating object proposals from edges. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014, Part V. LNCS*, vol. 8693, pp. 391–405. Springer, Heidelberg (2014)

Computer Vision – ECCV 2016

14th European Conference, Amsterdam, The

Netherlands, October 11–14, 2016, Proceedings, Part II

Leibe, B.; Matas, J.; Sebe, N.; Welling, M. (Eds.)

2016, XXIX, 887 p. 342 illus., Softcover

ISBN: 978-3-319-46474-9