

## Bottom-up Segmentation for Top-down Detection

Sanja Fidler<sup>1</sup>    Roozbeh Mottaghi<sup>2</sup>    Alan Yuille<sup>2</sup>    Raquel Urtasun<sup>1</sup>  
<sup>1</sup>TTI Chicago, <sup>2</sup>UCLA  
{fidler, rurtasun}@ttic.edu, {roozbehm@cs, yuille@stat}.ucla.edu

### Abstract

*In this paper we are interested in how semantic segmentation can help object detection. Towards this goal, we propose a novel deformable part-based model which exploits region-based segmentation algorithms that compute candidate object regions by bottom-up clustering followed by ranking of those regions. Our approach allows every detection hypothesis to select a segment (including void), and scores each box in the image using both the traditional HOG filters as well as a set of novel segmentation features. Thus our model “blends” between the detector and segmentation models. Since our features can be computed very efficiently given the segments, we maintain the same complexity as the original DPM [14]. We demonstrate the effectiveness of our approach in PASCAL VOC 2010, and show that when employing only a root filter our approach outperforms Dalal & Triggs detector [12] on **all** classes, achieving 13% higher average AP. When employing the parts, we outperform the original DPM [14] in 19 out of 20 classes, achieving an improvement of 8% AP. Furthermore, we outperform the previous state-of-the-art on VOC’10 test by 4%.*

### 1. Introduction

Over the past few years, we have witnessed a push towards holistic approaches that try to solve multiple recognition tasks jointly [29, 6, 20, 18, 33]. This is important as information from multiple sources should facilitate scene understanding as a whole. For example, knowing which objects are present in the scene should simplify segmentation and detection tasks. Similarly, if we can detect where an object is, segmentation should be easier as only figure-ground segmentation is necessary. Existing approaches typically take the output of a detector and refine the regions inside the boxes to produce image segmentations [22, 5, 1, 14]. An alternative approach is to use the candidate detections produced by state-of-the-art detectors as additional features for segmentation. This simple approach has proven very successful [6, 19] in standard benchmarks.

In contrast, in this paper we are interested in exploit-

ing semantic segmentation in order to improve object detection. While bottom-up segmentation has been often believed to be inferior to top-down object detectors due to its frequent over- and under- segmentation, recent approaches [8, 1] have shown impressive results in difficult datasets such as PASCAL VOC challenge. Here, we take advantage of region-based segmentation approaches [7], which compute a set of candidate object regions by bottom-up clustering, and produce a segmentation by ranking those regions using class specific rankers. Our goal is to make use of these candidate object segments to bias sliding window object detectors to agree with these regions. Importantly, unlike the aforementioned holistic approaches, we reason about all possible object bounding boxes (not just candidates) to not limit the expressiveness of our model.

Deformable part-based models (DPM) [14] and its variants [2, 35, 10], are arguably the leading technique to object detection<sup>1</sup>. However, so far, there has not been many attempts to incorporate segmentation into DPMs. In this paper we propose a novel deformable part-based model, which exploits region-based segmentation by allowing every detection hypothesis to select a segment (including void) from a small pool of segment candidates. Towards this goal, we derive simple features, which can capture the essential information encoded in the segments. Our detector scores each box in the image using both the traditional HOG filters as well as the set of novel segmentation features. Our model “blends” between the detector and the segmentation models by boosting object hypotheses on the segments. Furthermore, it can recover from segmentation mistakes by exploiting a powerful appearance model. Importantly, as given the segments we can compute our features very efficiently, our approach has the same computational complexity as the original DPM [14].

We demonstrate the effectiveness of our approach in PASCAL VOC 2010, and show that when employing only a root filter our approach outperforms Dalal & Triggs detector [12] by 13% AP, and when employing parts, we outperform the original DPM [14] by 8%. Furthermore, we outperform the previous state-of-the-art on VOC2010 by 4%.

<sup>1</sup>Poselets [4] can be shown to be very similar in spirit to DPMs

We believe that these results will encourage new research on bottom-up segmentation as well as hybrid segmentation-detection approaches, as our paper clearly demonstrates the importance of segmentation for object detection.

In the remainder of the paper, we first review related work and then introduce our novel deformable part-based model, which we call **segDPM**. We then show our experimental evaluation and conclude with future work.

## 2. Related Work

Deformable part-based model [14] and its variants have been proven to be very successful in difficult object detection benchmarks such as PASCAL VOC challenge. Several approaches have tried to augment the level of supervision in these models. Azizpour et al. [2] use part annotations to help clustering different poses as well as to model occlusions. Hierarchical versions of these models have also been proposed [35], where each part is composed of a set of sub-parts. The relative rigidity of DPMs has been alleviated in [10] by leveraging a dictionary of shape masks. This allows a better treatment of variable object shape. Desai et al. [13] proposed a structure prediction approach to perform non-maxima suppression in DPMs which exploits pairwise relationships between multi-class candidates. The tree structure of DPMs allows for fast inference but can suffer from problems such as double counting observations. To mitigate this, [27] consider lateral connections between high resolution parts.

In the past few years, a wide variety of segmentation algorithms that employ object detectors as top-down cues have been proposed. This is typically done in the form of unary features for an MRF [19], or as candidate bounding boxes for holistic MRFs [33, 21]. Complex features based on shape masks were exploited in [33] to parse the scene holistically in terms of the objects present in the scene, their spatial location as well as semantic segmentation. In [26], heads of cats and dogs are detected with a DPM, and segmentation is performed using a GrabCut-type method. By combining top-down shape information from DPM parts and bottom-up color and boundary cues, [32] tackle segmentation and detection task simultaneously and provide shape and depth ordering for the detected objects. Dai et al. [11] exploit a DPM to find a rough location for the object of interest and refine the detected bounding box according to occlusion boundaries and color information. [25] find silhouettes for objects by extending or refining DPM boxes corresponding to a reliably detectable part of an object.

DPMs provide object-specific cues, which can be exploited to learn object segmentations [3]. In [24], masks for detected objects are found by employing a group of segments corresponding to the foreground region. Other object detectors have been used in the literature to help segmenting object regions. For instance, while [4] finds segmentations

for people by aligning the masks obtained for each Poselet [4], [23] integrates low level segmentation cues with Poselets in a soft manner.

There are a few attempts to use segments/regions to improve object detection. Gu et al. [17] apply hough transform for a set of regions to cast votes on the location of the object. More recently, [28] learn object shape model from a set of contours and use the learned model of contours for detection. In contrast, in this paper we proposed a novel deformable-part based model, which allows each detection hypothesis to select candidate segments. Simple features express the fact that the detections should agree with the segments. Importantly, these features can be computed very efficiently, and thus our approach has the same computational complexity as DPM [14].

## 3. Semantic Segmentation for Object Detection

We are interested in utilizing semantic segmentation to help object detection. In particular, we take advantage of region-based segmentation approaches, which compute candidate object regions by bottom-up clustering, and rank those regions to estimate a score for each class. Towards this goal we frame detection as an inference problem, where each detection hypothesis can select a segment from a pool of candidates (those returned from the segmentation as well as void). We derive simple features, which can be computed very efficiently while capturing most information encoded in the segments. In the remainder of the section, we first discuss our novel DPM formulation. We then define our new segment-based features and discuss learning and inference in our model.

### 3.1. A Segmentation-Aware DPM

Following [14], let  $p_0$  be a random variable encoding the location and scale of a bounding box in an image pyramid as well as the mixture component id. As shown in [14] a mixture model is necessary in order to cope with variability in appearance as well as the different aspect ratios of the training examples. Let  $\{p_i\}_{i=1, \dots, P}$  be a set of parts which encode bounding boxes at double the resolution of the root. Denote with  $h$  the index over the set of candidate segments returned by the segmentation algorithm. We frame the detection problem as inference in a Markov Random Field (MRF), where each root filter hypothesis can select a segment from a pool of candidates. We thus write the score of a configuration as

$$E(\mathbf{p}, h) = \sum_{i=0}^P w_i^T \cdot \phi(x, p_i) + \sum_{i=1}^P w_{i,def}^T \cdot \phi(x, p_0, p_i) + w_{seg}^T \phi(x, h, p_0) \quad (1)$$

where  $h \in \{0, 1, \dots, H(x)\}$ , with  $H(x)$  the total number of segments for this class in image  $x$ . Note that  $h = 0$  im-

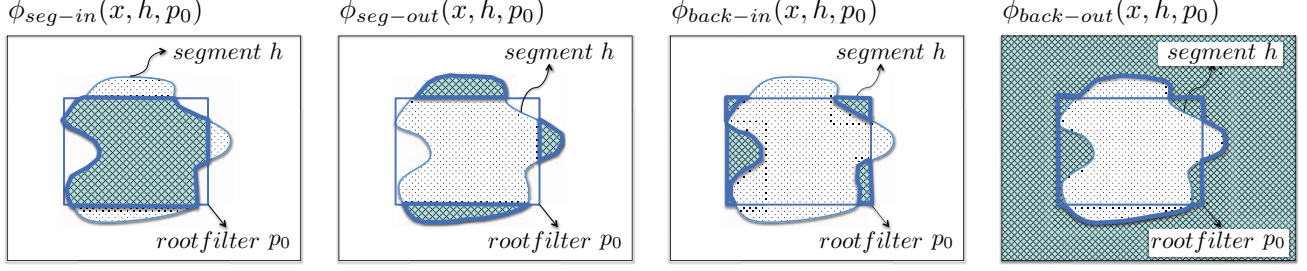


Figure 1. The box-segment features:  $\phi_{seg-in}$  and  $\phi_{seg-out}$ , encourage the box to contain as many segment pixels as possible. This pair of features alone could result in box hypotheses that “overshoot” the segment. The purpose of the second pair of features,  $\phi_{back-in}$  and  $\phi_{back-out}$ , is the opposite: it tries to minimize the number of background pixels inside the box and maximize its number outside. In synchrony these features would try to tightly place a box around the segment.

plies that no segment is selected. We will use  $S(h)$  to denote the segment that  $h$  indexes. As in [14], we employ a HOG pyramid to compute  $\phi(x, p_0)$ , and use double resolution to compute the part features  $\phi(x, p_i)$ . The features  $\phi(x, h, p_0)$  link segmentation and detection. In this paper, we define features at the level of the root, but our formulation can be easily extended to include features at the part level.

### 3.2. Segmentation Features

Given a set of candidate segments, we would like to encode features linking segmentation and detection while remaining computationally efficient. We would also like to be robust to over- and under- object segmentations, as well as false positive or missing segments. Towards this goal, we derive simple features which encourage the selected segment to agree with the object detection hypothesis. Most of our features employ integral images which makes them extremely efficient, as this computation can be done in constant time. We now describe the features in more details.

**Segment-In:** Given a segment  $S(h)$ , our first feature counts the percentage of pixels in  $S(h)$  that fall inside the bounding box defined by  $p_0$ . Thus

$$\phi_{seg-in}(x, h, p_0) = \frac{1}{|S(h)|} \sum_{p \in B(p_0)} \mathbb{1}\{p \in S(h)\}$$

where  $|S(h)|$  is the size of the segment indexed by  $h$ , and  $B(p_0)$  is the set of pixels contained in the bounding box defined by  $p_0$ . This feature encourages the bounding box to contain the segment.

**Segment-Out:** Our second feature counts the percentage of segment pixels that are outside the bounding box,

$$\phi_{seg-out}(x, h, p_0) = \frac{1}{|S(h)|} \sum_{p \notin B(p_0)} \mathbb{1}\{p \in S(h)\}$$

This feature discourages boxes that do not contain all segment pixels.

**Background-In:** We additionally compute a feature counting the amount of background inside the bounding box as follows

$$\phi_{back-in}(x, h, p_0) = \frac{1}{N - |S(h)|} \sum_{p \in B(p_0)} \mathbb{1}\{p \notin S(h)\}$$

with  $N$  the size of the image. This feature captures the statistics of how often the segments leak outside the true bounding box vs how often they are too small.

**Background-Out:** This feature counts the amount of background outside the bounding box

$$\phi_{back-out}(x, h, p_0) = \frac{1}{N - |S(h)|} \sum_{p \notin B(p_0)} \mathbb{1}\{p \notin S(h)\}$$

It tries to discourage bounding boxes that are too big and do not tightly fit the segments.

**Overlap:** This feature penalizes bounding boxes which do not overlap well with the segment. In particular, it computes the intersection over union between the candidate bounding box defined by  $p_0$  and the tighter bounding box around the segment  $S(h)$ . It is defined as follows

$$\phi_{overlap}(x, h, p_0) = \frac{B(p_0) \cap B(S(h))}{B(p_0) \cup B(S(h))} - \lambda$$

with  $B(S(h))$  the tighter bounding box around  $S(h)$ ,  $B(p_0)$  the bounding box defined by  $p_0$ , and  $\lambda$  a constant, which is the intersection over union level that defines a true positive. We employ in practice  $\lambda = 0.7$ .

**Background bias:** The value of all of the above features is 0 when  $h = 0$ . We incorporate an additional feature to learn the bias for the background segment ( $h = 0$ ). This puts the scores of the HOG filters and the segmentation potentials into a common referential. We thus simply define

$$\phi_{bias}(x, h, p_0) = \begin{cases} 1 & \text{if } h = 0 \\ 0 & \text{otherwise.} \end{cases}$$

Fig. 1 depicts our features computed for a specific bounding box  $p_0$  and segment  $S(h)$ . Note that the first two features,  $\phi_{seg-in}$  and  $\phi_{seg-out}$ , encourage the box to contain as many segment pixels as possible. This pair of features alone could result in box hypotheses that “overshoot” the segment. The purpose of the second pair of features,  $\phi_{back-in}$  and  $\phi_{back-out}$ , is the opposite: it tries to minimize the number of background pixels inside the box and maximize its number outside. In synchrony these features would try to tightly place a box around the segment. The overlap feature has a similar purpose, but helps us better tune the model to the VOC IOU evaluation setting.

### 3.3. Efficient Computation

Given the segments, all of our proposed features can be computed very efficiently. Note that the features have to be computed for each segment  $h$ , but this is not a problem as there are typically only a few segments per image. We start our discussion with the first four features, which can be computed in constant time using a single integral image per segment. This is both computationally and memory efficient. Let  $\phi_{int}(h)$  be the integral image for segment  $h$ , which, at each point  $(u, v)$ , counts the % of pixels that belong to this segment and are contained inside the subimage defined by the domain  $[0, u] \times [0, v]$ . This is illustrated in Fig. 2. Given the integral image  $\phi_{int}(h)$  for the  $h$ -segment, we compute the features as follows

$$\begin{aligned} \phi_{seg-in}(x, h, p_0) &= \phi_{br}(h, p_0) - \phi_{tr}(h, p_0) \\ &\quad - \phi_{bl}(h, p_0) + \phi_{tl}(h, p_0) \\ \phi_{seg-out}(x, h, p_0) &= |S(h)| - \phi_{seg-in}(x, h, p_0) \\ \phi_{back-in}(x, h, p_0) &= |B(p_0)| - \phi_{seg-in}(x, h, p_0) \\ \phi_{back-out}(x, h, p_0) &= (N - |S(h)|) - \phi_{back-in}(x, h, p_0) \end{aligned}$$

where as shown in Fig. 2,  $(\phi_{tl}, \phi_{tr}, \phi_{bl}, \phi_{br})$  indexes the integral image of segment  $S(h)$  at the four corners, i.e., top-left, top-right, bottom-left, bottom-right, of the bounding box defined by  $p_0$ .

The overlap feature between a hypothesis  $p_0$  and a segment  $S(h)$  can also be computed very efficiently. First, we compute the intersection as:

$$\begin{aligned} B(p_0) \cap B(S(h)) &= \\ \max[0, \min(x_{0,right}, x_{S(h),right}) - \max(x_{0,left}, x_{S(h),left})] \cdot \\ \max[0, \min(y_{0,bottom}, y_{S(h),bottom}) - \max(y_{0,top}, y_{S(h),top})] \end{aligned}$$

Note that the overlap will be non-zero only when each of the terms is larger than 0. Given that the segment bounding box  $B(S(h))$  is fixed and the width and height of  $p_0$  at a particular level of the pyramid are fixed as well, we can quickly

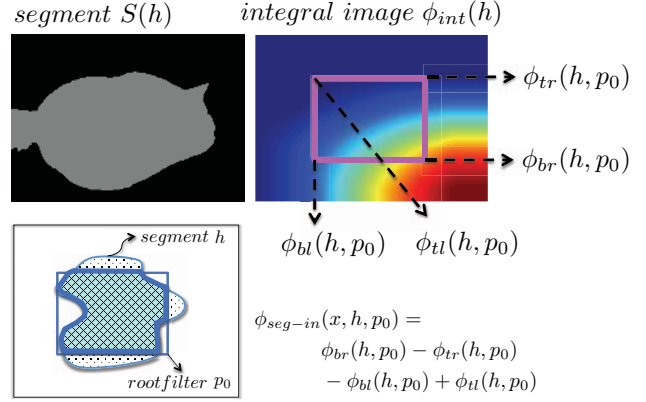


Figure 2. Segment feature computation via integral image.

compute the bounds of where in the image the feature needs to be computed (i.e., when the feature is different than 0). The denominator,  $B(p_0) \cup B(S(h))$ , can then be simply computed as the sum of the box areas minus the overlap.

### 3.4. Inference

Inference in our model can be done by solving the following optimization problem

$$\begin{aligned} \max_{p_0} \left( \sum_{i=0}^P w_i^T \cdot \phi(x, p_i) + \sum_{i=1}^P \max_{p_i} (w_{i,def}^T \cdot \phi(x, p_0, p_i)) + \right. \\ \left. + \max_h (w_{seg}^T \cdot \phi(x, h, p_0)) \right) \end{aligned}$$

Note that this can be done efficiently using dynamic programming as the structure of the graphical model forms a tree. The algorithm works as follows: First, we compute  $\max_h w_{seg}^T \phi(x, h, p_0)$  as well as  $\max_{p_i} (w_{i,def}^T \cdot \phi(x, p_0, p_i))$  for each root filter hypothesis  $p_0$ . We then compute the score as the sum of the HOG and segment score for each mixture component at the root level. Finally, we compute the maximum over the mixture components to get the score of an object hypothesis.

### 3.5. Learning

We learn a different weight for each feature using a latent structured-SVM [15]. Allowing different weights for the different segmentation features is important in order to learn how likely is for each class to have segments that undershoot or overshoot the detection bounding box. We employ as loss the intersection over the union of the root filters. As in DPM [14], we initialize the model by first training only the root filters, followed by training a root mixture model. Finally we add the parts and perform several additional iterations of stochastic gradient descent [14].

Note that we expect the weights for  $\phi_{seg-in}(x, h, p_0)$ ,  $\phi_{back-out}(x, h, p_0)$  and  $\phi_{overlap}(x, h, p_0)$  to be positive, as

we would like to maximize the overlap, the amount of foreground inside the bounding box and background outside the bounding box. Similarly, the weights for  $\phi_{seg-out}(x, h, p_0)$  and  $\phi_{back-in}(x, h, p_0)$  are expected to be negative as we would like to minimize the amount of background inside the bounding box as well as the amount of foreground segment outside. In practice, as the object’s shape can be far from rectangular, and the segments are noisy, the sign of the weights can vary to best capture the statistics of the data.

### 3.6. Implementation Details

We use CPMC [7] to get the candidate segments. In particular, for most experiments we use the final segmentation output of [7]. For each class, we find all connected components in the segmentation output, and remove those that do not exceed 1500 pixels. Unless otherwise noted, we do not use the score of the segments. On average, this gives us one segment per image. We also provide one experiment where we used more segments (5 on average per image), which we describe in Sec. 4.

## 4. Experimental Evaluation

We first evaluate our detection performance on `val` subset of PASCAL VOC 2010 detection dataset, and compare it to the baselines. We train all methods, including the baselines on the `train` subset. We use the standard PASCAL criterion for detection (50% IOU overlap) and report average precision (AP) as the measure of evaluation.

As baselines we use the Dalal&Triggs detector [12] (which for fairness we compare to our detector when only using the root filters), the DPM [14], as well as CPMC [7] when used as a detector. To compute the latter, we find all the connected components in the final segmentation output of CPMC [7], and place the tightest bounding box around each component. To compute the score of the box we utilize the CPMC ranking scores for the segments.

The comparison with [12] and our approach (segDPM) without parts is shown in the top Table 1, while the bottom table compares CPMC-based detector, DPM and our approach with parts. We significantly outperform the baselines: Dalal & Triggs detector by 13% and the CPMC baseline by 10%. Our model also achieves a significant boost of 8% AP over the DPM, which is a well established and difficult baseline to beat. Importantly, we outperform DPM in 19 out of 20 classes. The main power of our approach is that it blends between DPM (appearance) and segmentation (CPMC). When there is no segment, the method just scores a regular DPM. When there is a segment, our approach is encouraged to tightly fit a box around it. However, in cases of under- or over- segmentation, the appearance part of our model can still correctly position the box. Note that our results well demonstrate the effectiveness of using blended detection and segmentation models for object detection.

Fig. 4 depicts examples illustrating the performance of our approach. Note that our approach is able to both retrieve detections where there is no segment as well as position the bounding box correctly where there is segment evidence.

We evaluate our approach on VOC 2010 `test` in Table 2. Here, we trained CPMC [7], as well as our model on VOC `trainval`. We compare segDPM with DPM without the post-processing steps, i.e., bounding box prediction and context-rescoring, in the top of Table 2. In the bottom of Table 2 we compare our full approach with existing top scoring approaches. For the full approach, we show results when typical context-rescoring approach is used (same as in DPM), which we refer to as segDPM+rescore. We also show results when we rescored the detections by using the classification scores for each class, kindly provided to us by [9]. The classification (presence/absence of class in an image) accuracy measured by mean AP on VOC2010 is 76.2%. We refer to this approach with segDPM+rescore+classif. We outperform the competitors by 3.6%, and achieve the best result in 13 out of 20 classes.

We also experimented with using more segments, on the VOC 2010 `train / val` split. In particular, among 150 segments per image returned by [8], we selected a top-ranking subset for each class, so that there was an average of 5 segments per image. The results are reported in Table 3. We compare it to CPMC when using the same set of segments. One can see that with more segments our approach improves by 1.5%. As such, it outperforms DPM by 10%.

## 5. Conclusion

In this paper, we have proposed a novel deformable part-based model, which exploits region-based segmentation by allowing every detection hypothesis to select a segment (including void) from a pool of segment candidates. We derive simple yet very efficient features, which can capture the essential information encoded in the segments. Our detector scores each box in the image using both the HOG filters as in original DPM, as well as a set of novel segmentation features. This way, our model “blends” between the detector and the segmentation model, by boosting object hypotheses on the segments, while recovering from making mistakes by exploiting a powerful appearance model. We demonstrated the effectiveness of our approach in PASCAL VOC 2010, and show that when employing only a root filter our approach outperforms Dalal & Triggs detector [12] by 13% AP and when employing parts, we outperform the original DPM [14] by 8%. We believe that this is just the beginning of a new and exciting direction. We expect a new generation of object detectors which are able to exploit semantic segmentation yet to come.

**Acknowledgments** R.M. was supported in part by NSF 0917141 and ARL 62250-CS.

	plane	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	motor	person	plant	sheep	sofa	train	tv	Avg.
<b>VOC 2010 val, no parts</b>																					
Dalal [12]	29.1	36.9	2.9	3.4	15.6	47.1	27.1	11.4	9.8	5.8	6.0	5.0	24.8	28.4	27.5	2.2	18.4	9.2	27.4	23.2	18.1
segDPM (no parts)	<b>52.4</b>	<b>43.1</b>	<b>20.9</b>	<b>15.7</b>	<b>18.6</b>	<b>55.8</b>	<b>33.2</b>	<b>43.9</b>	<b>10.7</b>	<b>22.0</b>	<b>14.8</b>	<b>31.1</b>	<b>40.9</b>	<b>45.1</b>	<b>33.6</b>	<b>11.1</b>	<b>27.3</b>	<b>22.0</b>	<b>42.5</b>	<b>31.7</b>	<b>30.8</b>
<b>VOC 2010 val, with parts</b>																					
CPMC (no score) [7]	49.9	15.5	18.5	14.7	7.4	35.0	19.9	41.4	3.9	16.2	8.5	24.4	26.0	32.1	18.9	5.7	15.3	14.1	29.8	18.7	20.8
CPMC (score) [7]	53.3	19.5	22.8	15.7	8.1	42.7	22.1	<b>51.3</b>	4.3	18.9	10.5	28.1	30.5	38.3	20.9	6.0	19.2	18.6	35.4	21.1	24.4
DPM [14]	46.3	49.5	4.8	6.4	22.6	53.5	38.7	24.8	14.2	10.5	10.9	12.9	36.4	38.7	<b>42.6</b>	3.6	26.9	22.7	34.2	31.2	26.6
segDPM (parts)	<b>55.7</b>	<b>50</b>	<b>23.3</b>	<b>16.0</b>	<b>28.5</b>	<b>57.4</b>	<b>43.2</b>	49.3	<b>14.3</b>	<b>23.5</b>	<b>17.7</b>	<b>32.4</b>	<b>42.6</b>	<b>44.9</b>	42.1	<b>11.9</b>	<b>32.5</b>	<b>25.5</b>	<b>43.9</b>	<b>39.7</b>	<b>34.7</b>

Table 1. AP performance (in %) on VOC 2010 val for our detector with parts, the DPM [14], and the CPMC-based detector [7].

	plane	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	motor	person	plant	sheep	sofa	train	tv	Avg.
<b>VOC 2010 test, no post-processing</b>																					
DPM no postproc. [14]	47.2	<b>50.8</b>	8.6	12.2	<b>32.2</b>	48.9	44.4	28.1	13.6	22.7	11.3	17.4	40.4	47.7	<b>44.4</b>	7.6	30	17.3	38.5	34.3	29.9
segDPM no postproc.	<b>56.4</b>	48.0	<b>24.3</b>	<b>21.8</b>	31.3	<b>51.3</b>	<b>47.3</b>	<b>48.2</b>	<b>16.1</b>	<b>29.4</b>	<b>19.0</b>	<b>37.5</b>	<b>44.1</b>	<b>51.5</b>	<b>44.4</b>	<b>12.6</b>	<b>32.1</b>	<b>28.8</b>	<b>48.9</b>	<b>39.1</b>	<b>36.6</b>
<b>VOC 2010 test, with post-processing</b>																					
segDPM+rescore+classif	<b>61.4</b>	53.4	<b>25.6</b>	<b>25.2</b>	<b>35.5</b>	51.7	<b>50.6</b>	<b>50.8</b>	19.3	<b>33.8</b>	26.8	<b>40.4</b>	48.3	54.4	47.1	<b>14.8</b>	38.7	<b>35.0</b>	<b>52.8</b>	<b>43.1</b>	<b>40.4</b>
segDPM+rescore	<b>58.7</b>	51.4	<b>25.3</b>	<b>24.1</b>	33.8	52.5	<b>49.2</b>	<b>48.8</b>	11.7	30.4	21.6	37.7	46.0	53.1	46.0	13.1	35.7	29.4	<b>52.5</b>	41.8	<b>38.1</b>
NLPR_HOGLBP [34]	53.3	<b>55.3</b>	19.2	21.0	30.0	54.4	46.7	41.2	<b>20.0</b>	31.5	20.7	30.3	48.6	55.3	46.5	10.2	34.4	26.5	50.3	40.3	36.8
MITUCLA_HIERARCHY [35]	54.2	48.5	15.7	19.2	29.2	<b>55.5</b>	43.5	41.7	16.9	28.5	26.7	30.9	48.3	55.0	41.7	9.7	35.8	30.8	47.2	40.8	36.0
NUS_HOGLBP_CTX [9]	49.1	52.4	17.8	12.0	30.6	53.5	32.8	37.3	17.7	30.6	27.7	29.5	<b>51.9</b>	<b>56.3</b>	44.2	9.6	14.8	27.9	49.5	38.4	34.2
van de Sande et al. [30]	58.2	41.9	19.2	14.0	14.3	44.8	36.7	48.8	12.9	28.1	<b>28.7</b>	39.4	44.1	52.5	25.8	14.1	<b>38.8</b>	34.2	43.1	42.6	34.1
UOCTTL_SVM_MDPM [31]	52.4	54.3	13.0	15.6	35.1	54.2	49.1	31.8	15.5	26.2	13.5	21.5	45.4	51.6	<b>47.5</b>	9.1	35.1	19.4	46.6	38	33.7
Gu et al. [16]	53.7	42.9	18.1	16.5	23.5	48.1	42.1	45.4	6.7	23.4	27.7	35.2	40.7	49.0	32.0	11.6	34.6	28.7	43.3	39.2	33.1
UVA_DETMONKEY [30]	56.7	39.8	16.8	12.2	13.8	44.9	36.9	47.7	12.1	26.9	26.5	37.2	42.1	51.9	25.7	12.1	37.8	33.0	41.5	41.7	32.9
UVA_GROUPLOC [30]	58.4	39.6	18	13.3	11.1	46.4	37.8	43.9	10.3	27.5	20.8	36	39.4	48.5	22.9	13	36.8	30.5	41.2	41.9	31.9
BONN_FGT_SEG [8]	52.7	33.7	13.2	11.0	14.2	43.1	31.9	35.6	5.7	25.4	14.4	20.6	38.1	41.7	25.0	5.8	26.3	18.1	37.6	28.1	26.1

Table 2. AP performance (in %) on VOC 2010 test for our detector with parts and the DPM [14], without post processing (top table), and comparison with existing methods (only top 11 shown), with post-processing (table below).

## References

- [1] P. Arbelaez, B. Hariharan, C. Gu, S. Gupta, L. Bourdev, and J. Malik. Finding animals: Semantic segmentation using regions and parts. In *CVPR*, 2012. 1
- [2] H. Azizpour and I. Laptev. Object detection using strongly-supervised deformable part models. In *ECCV*, 2012. 1, 2
- [3] L. Bertelli, T. Yu, D. Vu, and B. Gokturk. Kernelized structural svm learning for supervised object segmentation. In *CVPR*, 2011. 2
- [4] L. Bourdev, S. Maji, T. Brox, and J. Malik. Detecting people using mutually consistent poselet activations. In *ECCV*, 2010. 1, 2
- [5] T. Brox, L. Bourdev, S. Maji, and J. Malik. Object segmentation by alignment of poselet activations to image contours. In *CVPR'11*. 1
- [6] G. Cardinal, X. Boix, J. van de Weijer, A. D. Bagdanov, J. Serrat, and J. Gonzalez. Harmony potentials for joint classification and segmentation. In *CVPR*, 2010. 1
- [7] J. Carreira, R. Caseiroa, J. Batista, and C. Sminchisescu. Semantic segmentation with second-order pooling. In *ECCV*, 2012. 1, 5, 6, 7
- [8] J. Carreira, F. Li, and C. Sminchisescu. Object Recognition by Sequential Figure-Ground Ranking. *IJCV*, 2011. 1, 5, 6
- [9] Q. Chen, Z. Song, Y. Hua, Z. Huang, and S. Yan. Hierarchical matching with side information for image classification. In *CVPR*, 2012. 5, 6
- [10] Y. Chen, L. Zhu, and A. Yuille. Active mask hierarchies for object detection. In *ECCV*, 2010. 1, 2
- [11] Q. Dai and D. Hoiem. Learning to localize detected objects. In *CVPR*, 2012. 2
- [12] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, pages I: 886–893, 2005. 1, 5, 6
- [13] C. Desai, D. Ramanan, and C. Fowlkes. Discriminative models for multi-class object layout. In *ICCV*, 2009. 2
- [14] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *PAMI*, 32(9), 2010. 1, 2, 3, 4, 5, 6
- [15] R. Girshick, P. Felzenszwalb, and D. McAllester. Object detection with grammar models. In *NIPS*, 2009. 4
- [16] C. Gu, P. Arbelaez, Y. Lin, K. Yu, and J. Malik. Multi-component models for object detection. In *ECCV*, 2012. 6
- [17] C. Gu, J. Lim, P. Arbelaez, and J. Malik. Recognition using regions. In *CVPR*, 2009. 2
- [18] G. Heitz, S. Gould, A. Saxena, and D. Koller. Cascaded classification models: Combining models for holistic scene understanding. In *NIPS*, 2008. 1
- [19] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NIPS*, 2011. 1, 2
- [20] L. Ladicky, C. Russell, P. Kohli, and P. H. Torr. Graph cut based inference with co-occurrence statistics. In *ECCV*, 2010. 1
- [21] L. Ladicky, P. Sturgess, K. Alahari, C. Russell, and P. H. Torr. What, where and how many? combining object detectors and crfs. In *ECCV*, 2010. 2
- [22] V. Lempitsky, P. Kohli, C. Rother, and B. Sharp. Image segmentation with a bounding box prior. In *ICCV*, 2009. 1
- [23] M. Maire, S. X. Yu, and P. Perona. Object detection and segmentation from joint embedding of parts and pixels. In *ICCV*, 2011. 2
- [24] A. Monroy and B. Ommer. Beyond bounding-boxes: Learning object shape by model-driven grouping. In *ECCV12*. 2
- [25] R. Mottaghi. Augmenting deformable part models with irregular-shaped object patches. In *CVPR*, 2012. 2
- [26] O. Parkhi, A. Vedaldi, C. V. Jawahar, and A. Zisserman. The truth about cats and dogs. In *ICCV*, 2011. 2
- [27] M. Pedersoli, A. Vedaldi, and J. Gonzalez. A coarse-to-fine approach for fast deformable object detection. In *CVPR*, 2011. 2
- [28] P. Srinivasan, Q. Zhu, and J. Shi. Many-to-one contour matching for describing and discriminating object shape. In *CVPR*, 2010. 2
- [29] E. Sudderth, A. Torralba, W. T. Freeman, and A. Wilsky. Learning hierarchical models of scenes, objects, and parts. In *ICCV*, 2005. 1

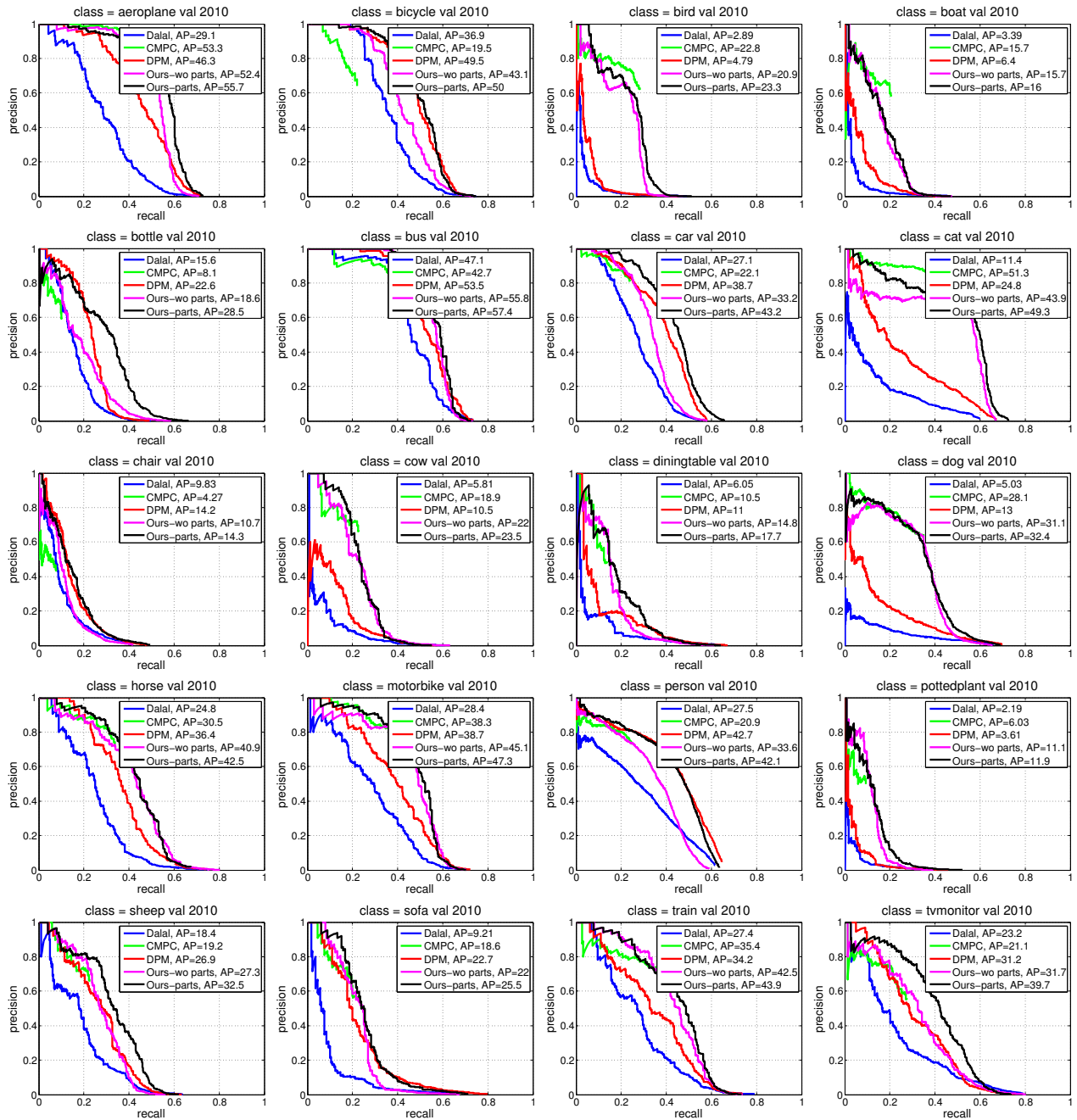


Figure 3. Precision-recall curves on PASCAL VOC 2010 val. Note that our approach significantly outperforms all baselines.

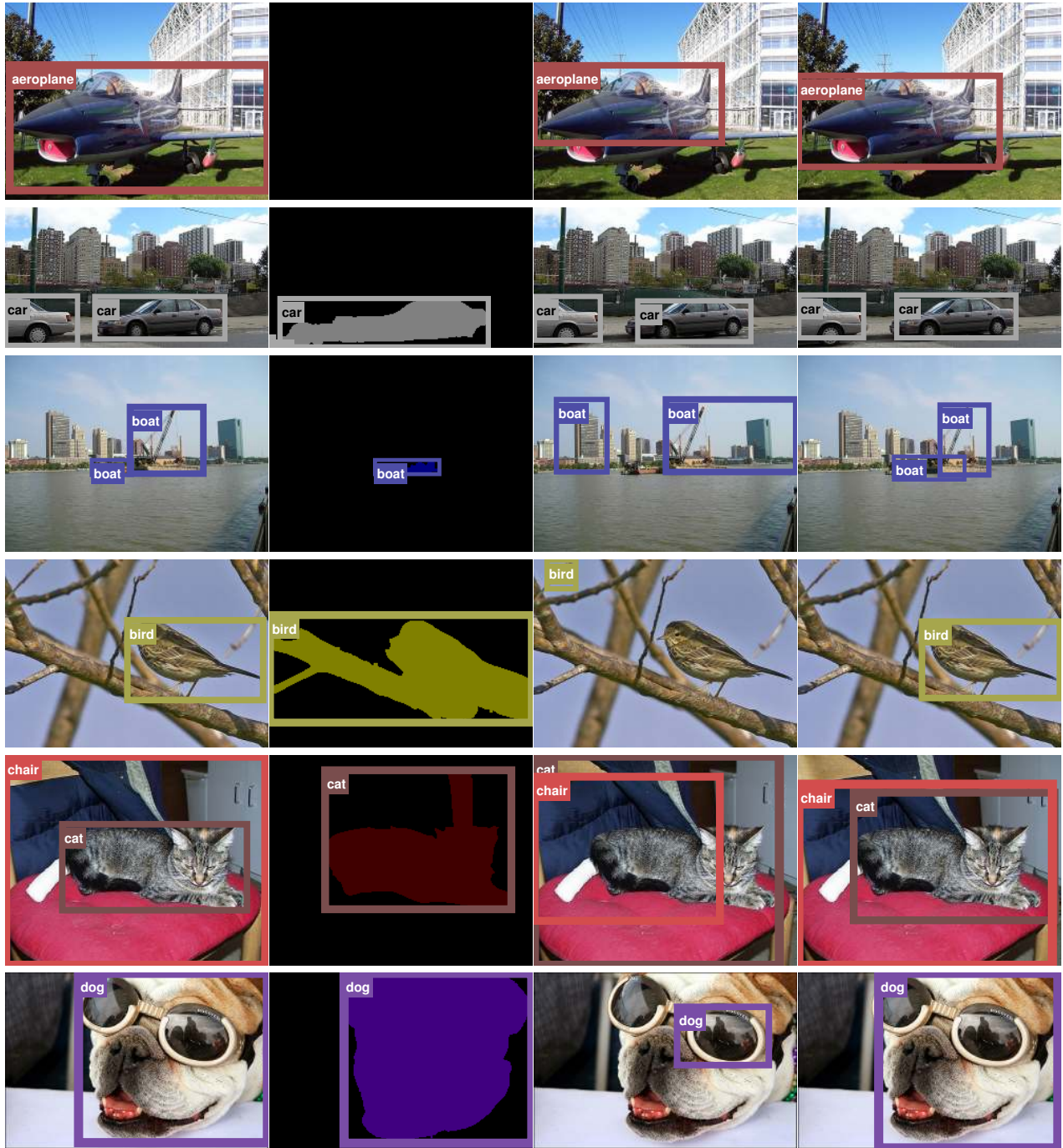
	plane	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	motor	person	plant	sheep	sofa	train	tv	Avg.
<b>VOC 2012 val, more segments</b>																					
CPMC (1 seg) [7]	53.3	19.5	22.8	15.7	8.1	42.7	22.1	51.3	4.3	18.9	10.5	28.1	30.5	38.3	20.9	6.0	19.2	18.6	35.4	21.1	24.4
CPMC (5 seg) [7]	59.8	27.6	27.1	19.6	12.7	53.1	31.2	56.6	8.2	25.6	17.5	34.8	39.8	42.3	25.9	10.3	29.8	26.6	46.7	33.4	31.4
segDPM (1 seg)	55.7	50.0	23.3	16.0	28.5	57.4	43.2	49.3	14.3	23.5	17.7	32.4	42.6	44.9	42.1	11.9	32.5	25.5	43.9	39.7	34.7
segDPM (5 seg)	56.1	49.0	22.9	18.2	34.0	58.9	42.9	49.8	15.4	25.0	22.7	32.3	46.2	45.6	39.2	13.6	33.3	30.6	46.7	41.5	<b>36.2</b>

Table 3. AP performance (in %) on **VOC 2010 val** for our detector when using more segments.

[30] K. E. A. van de Sande, J. R. R. Uijlings, T. Gevers, and A. W. M. Smeulders. Segmentation as selective search for object recognition.

In *ICCV*, 2011. 6

[31] A. Vedaldi, V. Gulshan, M. Varma, and A. Zisserman. Multiple ker-



(a) GT (b) CPMC (c) DPM (d) segDPM

Figure 4. For each method, we show top  $k$  detections for each class, where  $k$  is the number of boxes for that class in GT. For example, for an image with a chair and a cat GT box, we show the top scoring box for chair and the top scoring box for cat.

nels for object detection. In *ICCV*, 2009. 6

[32] Y. Yang, S. Hallman, D. Ramanan, and C. Fowlkes. Layered object models for image segmentation. *PAMI*, 2011. 2

[33] J. Yao, S. Fidler, and R. Urtasun. Describing the scene as a whole: Joint object detection, scene classification and semantic segmentation. In *CVPR*, 2012. 1, 2

[34] Y. Yu, J. Zhang, Y. Huang, S. Zheng, W. Ren, C. Wang, K. Huang, and T. Tan. Object detection by context and boosted hog-lbp. In *ECCV w. on PASCAL*, 2010. 6

[35] L. Zhu, Y. Chen, A. Yuille, and W. Freeman. Latent hierarchical structural learning for object detection. In *CVPR*, 2010. 1, 2, 6