

Contour-Based Object Detection as Dominant Set Computation

Xingwei Yang¹, Hairong Liu², Longin Jan Latecki¹

¹Department of Computer and Information Sciences, Temple University,
Philadelphia, PA, 19122.

Email: {xingwei,latecki}@temple.edu

²Department of Electrical and Computer Engineering, National University of
Singapore, Singapore

Email: lhrbss@gmail.com

Abstract. Contour-based object detection can be formulated as a matching problem between model contour parts and image edge fragments. We propose a novel solution by treating this problem as the problem of finding dominant sets in weighted graphs. The nodes of the graph are pairs composed of model contour parts and image edge fragments, and the weights between nodes are based on shape similarity. Because of high consistency between correct correspondences, the correct matching corresponds to a dominant set of the graph. Consequently, when a dominant set is determined, it provides a selection of correct correspondences. As the proposed method is able to get all the dominant sets, we can detect multiple objects in a image in one pass. Moreover, since our approach is purely based on shape, we also determine an optimal scale of target object without a common enumeration of all possible scales. Both theoretic analysis and extensive experimental evaluation illustrate the benefits of our approach.

1 Introduction

Object detection in cluttered images, with scale and intra-class variations, is one of the most difficult problems in computer vision. Appearance based methods have had remarkable success in recent years [1–5]. However, in many cases, the appearance between intra-class objects varies a lot [6], which makes the appearance features not reliable. Thus, recently we have observed a significant increase in methods that utilize contour shape [7–11]. However, shape based methods also face many challenges, such as pose variance, missing edges, and view point changes. Among these challenges, a critical one seems to be missing contour fragments in the cluttered edge images. The contour fragments may be missing due to occlusion or due to missing edges, since important contours of target objects are often hard or impossible to detect by state-of-the-art edge detectors [12].

Interestingly, our visual system can perform contour grouping, object detection, and recognition, even if only cluttered edge information is provided, e.g., Fig. 1. We can easily perform all these tasks even if the important contour information is missing, and we may not be able to complete it, e.g., we can recognize

the giraffes in Fig. 1, but we may not be able to draw or imagine the missing outline of their heads. Thus, we can perform contour grouping, object detection, and recognition while keeping at least part of missing information ambiguous, and we do not attempt to disambiguate all missing information. In other words, we do not attempt to completely reconstruct all contour parts of the object in the image. This fact is one of the key motivations for the proposed inference method.



Fig. 1. Parts of the object contours are missing due to missing edges.

We formulate object detection as a labeling or matching problem between image segments and model parts. As we discussed, in order to achieve a human like performance in recognizing objects, it requires computing partial assignment between the image segments and model parts, which has been a critical problem for traditional labeling methods [13–15]. To deal with this problem, we propose to transform the matching problem into finding dominant sets in a correspondence graph, in which each vertex represents a pair of image and model segments and the affinities between vertices are obtained by shape similarity. With this modification, missing parts of a true object contour in the image do not negatively influence the selection of dominant sets. The concept of dominant sets has been introduced in [16], where also a method for dominant set computation is proposed. Each dominant set is a local solution of a constrained quadratic function, and it is computed by a recursive procedure that depends on the initialization. Recently [17] proposed a novel initialization strategy that is guaranteed to yield all dominant sets under certain assumptions. Although the assumptions in [17] are derived for the application of common visual pattern discovery, they also apply to our application. Different from [16] and from [17], where dominant sets are treated as final solutions, we view each dominant set as a solution hypothesis that is evaluated with global shape similarity. The main reason is that the value of the target quadratic function is based only on local shape similarity, which may be insufficient for object detection. In other words, the dominant set with the highest value of the target quadratic function does not necessarily imply a correct object detection. This is also the reason why we need to consider all dominant sets.

There are at least three key advantages of the proposed method. It is insensitive to noise and outliers, thus, it can detect objects in cluttered images. The fact that we consider all dominant sets provides another main advantage. It allows us to detect multiple objects in one pass, which is also difficult for traditional la-

belonging methods. Each object instance is represented by a different dominant set. Moreover, as the proposed method is purely based on shape similarity, we can automatically detect objects in different scales without enumerating the scales, which deals with the problem of resolution. It is an important benefit compared to other methods, such as sliding window [18] and Hough voting [19]. Both methods have to explicitly enumerate various scales in a certain scale range to obtain the best solution. They require a predetermined scale range, which is a hidden parameter not mentioned in most papers.

Example images in Fig. 2 demonstrate the benefit of the proposed method. The white lines are edge segments after edge linking (see Section 3) and the red lines are the detected segments. In these images, parts of the objects are missing due to occlusion, and some of them contain multiple objects like the four apple logos in the top left. Our approach can detect multiple objects at the same time, i.e., in one pass, and allow partial assignment between image segments and model parts, which makes the system robust to missing edges and occlusion. For example, the two mugs are partially occluded, and many bottles are missing some edges.

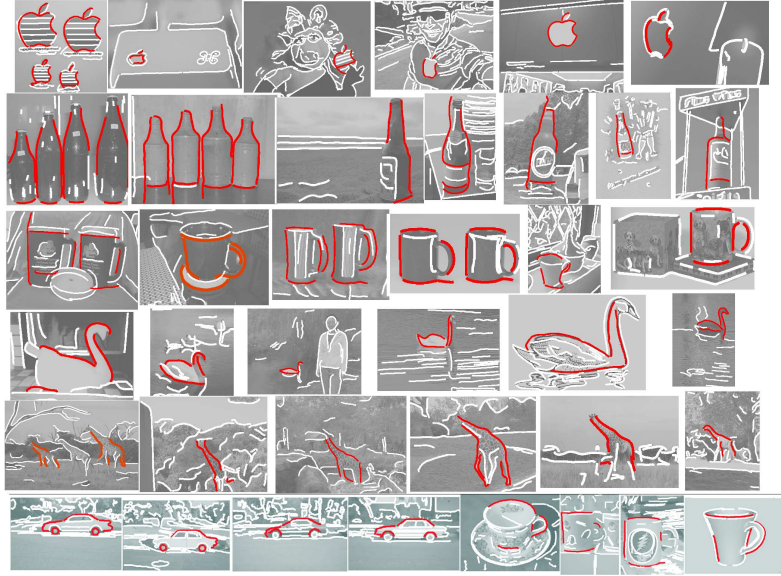


Fig. 2. Example detections on ETHZ dataset and, in the last row, on two Caltech-101 classes: car-side and cups.

The paper is organized as follow. Related methods are discussed in Section 2. The pre-processing step is described in Section 3. We introduce our approach based on dominant sets in Section 4. The optimization method for finding the dominant sets is described in Section 5. We view the dominant sets as object detection hypotheses, which are then evaluated with global shape similarity in

Section 6. In Section 7, we evaluate the performance of the proposed approach on the challenging ETHZ shape dataset [20, 7], which features large variations in scale and cluttered background. Besides, we also evaluate our approach on a subset of Caltech 101 dataset [21].

2 Related Work

As there exists a lot of papers on shape based object detection and recognition, we only review the most related ones. Ferrari et al. [20] propose to use kAS, the k connected roughly straight contour segments, with Hough voting to detect objects. Later, Ferrari et al. [7] extend their work to learn the model from the image. To improve [7], Jiang et al. [10] propose to learn a shape prior model for each object class. Boundary fragments combined with classifier have also been investigated in [22]. Instead of object's contour, Trinh and Kimia [23] use skeleton-based generative shape model with modified dynamic programming to detect objects. Bai. et al. [24] also utilize skeleton to constrain the detection process. All the above methods require multiple initializations and they enumerate all possible object sizes (scales) to get the optimal results. Different from them, the proposed method is able to detect multiple objects at different scales in one pass without enumerating scales.

Ravishankar et al. [25] introduce a multi-stage contour based detection approach with dynamic programming, which is also scale independent. Different from them, we solve the matching problem by finding dominant sets in the correspondence graph. Zhu et al. [8] utilize Shape Context [26] to evaluate the distance between model and image segments. They formulate the shape matching of contours as a set-set matching problem and solve it by linear programming, which is fundamentally different from us. Similar to our method, Lu et al. [11] formulate object detection as a segment correspondence problem. However, their inference framework is very different, where they utilize particle filter to solve the label assignment problem. Furthermore, they cannot detect multiple objects.

Gu et al. [27] utilize region segmentation to detect target objects. An appearance based approach was recently used by Maji and Malik [28] by integrating Hough transform based features of codebooks into kernel classifiers. To solve the problem of scales in Hough voting, Ommer and Malik [29] propose a weighted, pairwise clustering of voting lines to obtain globally consistent hypotheses. Then, a verification stage is used to re-rank the hypotheses. Unlike [29, 28, 27], we use a purely shape based method and do not utilize any classifiers like SVM to rank the hypotheses.

3 Preprocessing

As we formulate the object detection as a correspondence problem between image segments and model segments, we need to construct image segments from image edge maps as well as define shape models composed of contour segments. We utilize shape similarity of these segments to perform object detection. Given an

image I and the edge map, we use an open source edge linking method provided by Peter Kovess [30] to group edge pixels into edge fragments. If a junction point exists on the edge fragment, the corresponding edge fragments are split at the junction point. We obtain a set of image edge segments $E = \{e_1, \dots, e_n\}$ for the image I . An example is shown in Fig. 3(a), where each edge segment is shown in a different color.

The model segments $S = \{s_1, \dots, s_m\}$ are manually designed so that they represent meaningful contour parts. As junction points normally exist at high curvature points in the edge maps, we also decompose the model template at high curvature points. Moreover, since the image segments are noisy and some part of object boundary may be missing, we need to add shorter model segments in addition to longer ones. Then, the segments are grouped into different part bundles $\mathcal{B} = \{B_k\}_{k=1}^b$, where each part bundle represents the same visual part of the modeled contour. The main constraint for the bundle design is to ensure a rough shape sketch constructed by selecting one part from each bundle can still resemble the model contour. An example is shown in Fig. 3(b).

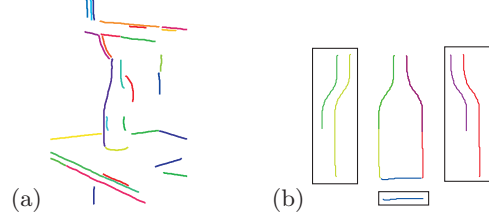


Fig. 3. (a) The edge segments in one image. Different colors represent different segments; (b) Model segments for category bottle, which is shown in the middle. Each segment is shown in a different color and the segments in one box form a part bundle. Thus, there are three part bundles for the bottle.

4 Problem Formulation

With preprocessing introduced in Section 3, we can obtain a set of image segments $E = \{e_1, \dots, e_n\}$ for the image I and a set of model segments $S = \{s_1, \dots, s_m\}$ for a model contour template. We formulate the object detection as a labeling problem, labeling the model segments with the image segments. Our goal is to find the segment correspondence so that the image segment corresponding to the model segments maximize the global shape similarity of all selected model segments to all selected image segment.

To reach this goal, we first build a graph G whose nodes are $\{c_1, \dots, c_m\}$, where $c_i = (s_i, e_{i'}) \in S \times E$, and $i = 1, \dots, m$. A tangent distance $TD(s_i, e_{i'})$ between two segments s_i and $e_{i'}$ is computed by matching sequences of their tangent directions with dynamic programming. Since we resample all the segments to the same number of sample points, TD is scale invariant. For each

model segment s_i , we use TD to find K most similar image segments. Hence graph G has $M = m \times K$ nodes and each node represents a correspondence c_k for $k = 1, \dots, M$. The weight of an edge connecting nodes i and j is defined as:

$$w_{ij} = \mathcal{N}(SC(s_i \cup s_j, e_{i'} \cup e_{j'})) \quad (1)$$

where $c_i = (s_i, e_{i'})$ and $c_j = (s_j, e_{j'})$ are two correspondences, \mathcal{N} is a Gaussian, and SC is the Shape Context [26] distance. The mean of Gaussian \mathcal{N} is defined as 0 and the standard deviation is defined as a quarter of the average distance between all pairs of correspondences.

Hence w_{ij} represents shape similarity between shape constructed of two model segments $s_i \cup s_j$ and two image edge segments $e_{i'} \cup e_{j'}$. The adoption of Shape Context have several advantages. First, it is a descriptive shape similarity method and it can be easily used for discrete points sets. Second, SC performs automatic scale normalization. Consequently, w_{ij} is scale invariant.

However, not all correspondences in graph G are compatible. For example, two edge segments that are far away from each other in a given image cannot both belong to the contour of a target object whose diameter is smaller than their distance. Therefore, we will define now a binary relation that allows us to efficiently remove such correspondences from G . We observe that when computing the shape distance $TD(s_i, e_{i'})$, we also obtain the correspondence of sample points of s_i to sample points of $e_{i'}$. It allows us to determine the scale factor so that the model bounding box can be properly re-scaled. Then the re-scaled bounding box is placed on the image; since we know the position of the model bounding box relative to segment s_i , the position of the bounding box in the image is determined by the position of $e_{i'}$.

Let us denote two re-scaled and relocated model bounding boxes in the image with $bbx(i)$ and $bbx(j)$, which are based on $c_i = (s_i, e_{i'})$ and $c_j = (s_j, e_{j'})$, respectively. Of course, if both correspondences are correct the two bounding boxes in the image should coincide. In the left image of Fig. 4, the color segments are the image segments being considered and the numbers show the indices of the corresponding model segments, which are 1 and 3. The two estimated bounding boxes are shown in red and green. Although the bounding box estimation is not perfect, it can roughly determine that the two image edge segments can belong to the same contour. On the other hand, the right image shows a wrong correspondence that leads to two disjoint bounding boxes. Therefore, if the area of the intersection of both bounding boxes in the image is small, then $e_{i'}$ and $e_{j'}$ cannot be both parts of the contour of the target object. To capture this property, we define a binary relation

$$R_I(i, j) = \begin{cases} 1, & \frac{area(bb x(i) \cap bb x(j))}{area(bb x(i) \cup bb x(j))} > C \\ 0, & otherwise, \end{cases} \quad (2)$$

where C is the area intersection threshold that is set to 0.1 in all our experiments.

We also define another binary relation that relates model segments of two correspondences. Since the shape constructed by two segments $s_i \cup s_j$ that belong

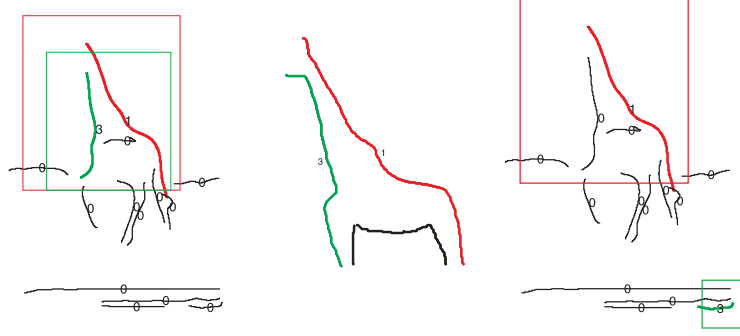


Fig. 4. Left: The estimated bounding boxes of two correct correspondences; Right: The estimated bounding boxes of two wrong correspondences. Middle: The shape model with corresponding segments marked in colors.

to the same part bundle B_k for $k = 1, \dots, b$ is not particularly informative, because they represent the same model part, we do not allow correspondences that involve such model segments. We define $R_M(i, j) = 0$ if $s_i, s_j \in B_k$ for some $k = 1, \dots, b$ and $R_M(i, j) = 1$ otherwise.

The weighted adjacency matrix A of graph G is an $M \times M$ matrix defined as

$$A_{i,j} = \begin{cases} 0, & i = j \text{ or } R_I(i, j) = 0 \text{ or } R_M(i, j) = 0 \\ w_{ij}, & \text{otherwise.} \end{cases} \quad (3)$$

The binary relations R_I and R_M help us to make the graph G sparse, which significantly reduces the computation cost. It is obvious that A is symmetric and nonnegative, since this is the case for w_{ij} , R_I and R_M .

We are interested in finding subgraphs H of G that are local maxima of the average affinity score S_a defined as

$$S_a(H) = \frac{1}{|H|^2} \sum_{i \in H, j \in H} A_{ij} = \mathbf{x}^T A \mathbf{x}, \quad (4)$$

where $|H|$ is the number of nodes of H and \mathbf{x} is a column vector such that $x_i = 1/|H|$ if $i \in H$ and $x_i = 0$ otherwise for $i = 1, \dots, M$.

For unweighted graphs, the Motzkin-Straus theorem [31] has established a connection between the maximal cliques and the local maximizers of the quadratic function:

$$\text{maximize } f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} \quad \text{subject to } \mathbf{x} \in \Delta, \quad (5)$$

where $\Delta = \{\mathbf{x} \in \mathbb{R}^m : \mathbf{x} \geq 0 \text{ and } |\mathbf{x}|_1 = 1\}$ is the standard simplex in \mathbb{R}^m . Eq. 5 means that a subgraph H of G is a maximal clique if and only if its characteristic vector \mathbf{x}^H is a local maximizer of this equation, where $x_i^H = 1/|H|$ if $i \in H$, and $x_i^H = 0$ otherwise. Recently, Pavan and Pelillo [16] generalized the

Motzkin-Straus theorem to weighted graphs. They also introduced the notation of dominant sets of vertices as a generalization to weighted graphs of the concept of a maximal cliques in unweighted graphs. In unweighted graphs dominant sets are equivalent to (strictly) maximal cliques. They showed that each (strict local) solution of the quadratic program Eq. 5 determines a **dominant set**, which we take as a definition of the dominant set in this paper.

As has been observed in [16, 17], Eq. 5 maximizes the same quadratic function as the spectral methods in [32, 33]. The only difference is the constraints on \mathbf{x} : the spectral methods require $\|\mathbf{x}\|_2 = 1$ while Eq. 5 requires $\|\mathbf{x}\|_1 = 1$. This minor difference changes dramatically the properties of obtained solutions. Instead of partitioning all data, as is the case for spectral methods, Eq. 5 only selects highly correlated data and ignores outliers. Consequently, the proposed object detection system can automatically select contours of the target object in an edge image and at the same time ignore the vast majority of the background segments.

5 Optimization

The main challenge we face now is to determine all local maxima of the quadratic program Eq. 5. Following [16], once an initialization $\mathbf{x}(1)$ is given at discrete time step 1, the discrete replicator equation [34] can be used to obtain a local solution \mathbf{x}^* :

$$\mathbf{x}_i(t+1) = \mathbf{x}_i(t) \frac{(A \mathbf{x}(t))_i}{\mathbf{x}(t)^T A \mathbf{x}(t)} \quad (6)$$

for $i = 1, \dots, M$ indexing the coordinates of vector \mathbf{x} . As is proven in [16], each strict local solution of (Eq. 5) determines a dominant set.

A key question for our approach is how to enumerate the initialization vectors $\mathbf{x}(1)$ so that we can obtain all local maxima $\{\mathbf{x}^*\}$. To solve this problem, Pavan and Pelillo [16] propose to detect dominant sets iteratively, i.e., after finding a dominant set, they remove its vertices from the graph G , and then rerun the algorithm on the remaining vertices. However, their method may miss some local maxima due to the fact that some dominant sets have nonempty intersection.

Recently [17] proposed a novel initialization strategy, which we follow in our approach. They suggest initializing $\mathbf{x}(1)$ in the neighborhood $N(v) \cup \{v\}$ of every vertex $v \in G$, where $N(v)$, the neighborhood of v is determined by thresholding the row v of matrix A . This strategy is based on the assumption that each dominant set that contains v must be a subset of $N(v) \cup \{v\}$. This assumption is satisfied in our setting for dominant sets that correctly determine a target shape composed of edge segments in the image, since all correspondences that extend the correspondence v to form the target shape have relatively large affinity with v . Consequently, this initialization strategy does not eliminate any correct solution in our setting.

However, this initialization strategy produces many duplicate solutions. Therefore, [17] proposed merging two dominant sets if their indicator vectors \mathbf{x}^* and \mathbf{y}^* have large correlation defined as $(\mathbf{x}^*)^T \mathbf{y}^*$. This reflects the fact that although different dominant sets may have nonempty intersection, their overlap is usually

small. This merging strategy significantly reduces the set of solution hypotheses that we need to examine with global shape similarity as described in the next section.

6 Final Evaluation

All different dominant sets obtained as solutions to Eq. 5 are potential object detection hypotheses in our system. Our final step is to use global shape similarity to evaluate these hypotheses. Although computing the global shape similarity is computationally expensive, it is tractable in our application. Usually the graph of correspondences G may have several hundred vertices, but we obtain less than 20 different dominant sets as solutions to Eq. 5.

Since each coordinate \mathbf{x}_i^* of \mathbf{x}^* represents the probability that correspondence i has been selected, we simply select the correspondences with probability bigger than 0 as the elements of the dominant set L determined by \mathbf{x}^* . We obtain $L = \{c_{i_1}, \dots, c_{i_l}\}$ for some $l \ll M$, where $c_{i_k} = (s_{i_k}, e_{i'_k})$. The global shape distance is computed with shape context as

$$SC(L) = SC(\bigcup_{k=1}^l s_{i_k}, \bigcup_{k=1}^l e_{i'_k}). \quad (7)$$

Thus, we simply compare the shape formed by combining all selected image edge segments $\bigcup_{k=1}^l e_{i'_k}$ to the corresponding model segments $\bigcup_{k=1}^l s_{i_k}$.

Although in Eq. 3, we have removed the connection between the model segments from the same part bundle, L can still contain different model segments from the same part bundle. Since each part bundle can only provide at most one part for one detection, we need to repeat the final evaluation Eq. 7 for each model segment from the same part bundle separately, and then select the best score. Moreover, the more different part bundles appear in the correspondences in L , the more reliable is the global shape similarity. Thus, we only consider dominant sets L that contain correspondences from a certain minimal number of part bundles, which is set to 3 in all our experimental results. The global shape distance in Eq. 7 is used to rank the object detection hypotheses. We stress that our method is purely based on shape similarity without using any training or classifier.

7 Experimental Results

To evaluate our approach, we choose the challenging ETHZ Shape Dataset [7, 20] containing five different categories with 255 images in total. Each image contains one or more instances with significant background clutter. All categories have significant scale difference and intra-class variation. Similar to Zhu et al. [8] and Lu et al. [11], we use a single manually constructed, contour model for each shape class. The detection performance is measured based on the standard PASCAL VOC criterion [35].

Furthermore, we also selected two classes from Caltech-101 dataset [21] to evaluate our approach: cups and car-sides. They contain substantial intra-class variations and missing edge segments. Similar to [20] an equal number of negative images is selected from the Caltech-101 background set. Then, the test set for each class consists of all positive images and the equal number of negative images. The contour model for cups is the same as the model for ETHZ mugs, and we manually created a contour model for car-sides.

7.1 ETHZ Dataset

In Fig. 2, some example detections on ETHZ dataset are shown. These examples demonstrate the proposed method can handle multiple detections and scale variation in one pass. Moreover, our method is robust to missing segments, even if a whole meaningful object part is missing, which cannot be solved by relaxation labeling and many other methods.

We also demonstrate the benefit of our algorithm by comparing to other methods. A large number of methods have been tested on ETHZ dataset [27, 28, 11, 7, 8, 20, 29]. It is difficult to compare to all of them, thus, we only compare to some shape based methods. We first compare to [20, 7, 29] by plotting the detection rate (DR) against false positive per image (FPPI), see Fig. 5. Besides the curves, the detection rates at 0.3 and 0.4 FPPI are also shown in Table 1. We stress that only the method by [29] and the proposed method are truly scale independent. The other two methods [20, 7] enumerate scales in a certain scale range.

Table 1. Detection rates at 0.3/0.4 FPPI for ETHZ dataset. The best results are highlighted in bold.

Category	Clustering Lines [29]	KAS [20]	Full system [7]	Our method
Apples	95.0/95.0	50.0/60.0	77.7/83.2	80.0/80.0
Bottles	89.3/89.3	92.9/92.9	79.8/81.6	92.9/95.9
Giraffes	70.5/75.4	49.0/51.1	39.9/44.5	76.92/79.21
Mugs	87.3/90.3	67.8/77.4	75.1/80.0	83.3/84.85
Swans	94.1/94.1	47.1/52.4	63.2/70.5	90.9/ 94.1
Average	87.2/88.8	61.4/66.8	67.2/72.0	84.8/86.79

Ferrari et al. [20, 7] train their detector for each category on half of the positive examples on that class. In [20], they also use the negative images for training. In comparison to [7], it turns out that our algorithm can obtain better results on the whole ETHZ dataset except the class applelogos, where we perform equally well. Our method improves the average detection rate by 17.7% and 14.9% at 0.3 and 0.4 FPPI, respectively. Similarly, we are comparable to [20] on the class bottle and better on all the other classes. The average detection rate has increased by 23.5% and 20.1% at 0.3 and 0.4 FPPI respectively.

Furthermore, we also compared to the recent work by [29], which also uses half of the positive examples as training. Our method performs better on classes

giraffes and bottles and is comparable on class swans. However, we perform worse on mugs and apples. The average detection rates at 0.3 and 0.4 FPPI are about 2% lower compared to [29].

We stress that our algorithm is purely shape based and we do not have any postprocessing phase to refine the results. In [29], the SVM classifier is run in sliding window mode over a grid of locations around each initial detection, which boosts the results a lot. To make the comparison to [29] complete, we also compare the precision/recall curves in Fig. 6. Unlike the previous comparison in DR/FPPI, the actual values for precision/recall are not reported in [29]. Thus, we can only compare the curves visually. Our results on swans and bottles are better than [29] and the precision of giraffes is better with a little lower recall value. The precision of apple logos is comparable, but the recall is worse. Our results of mugs are worse than [29]. We also compared to methods in [8, 11] using

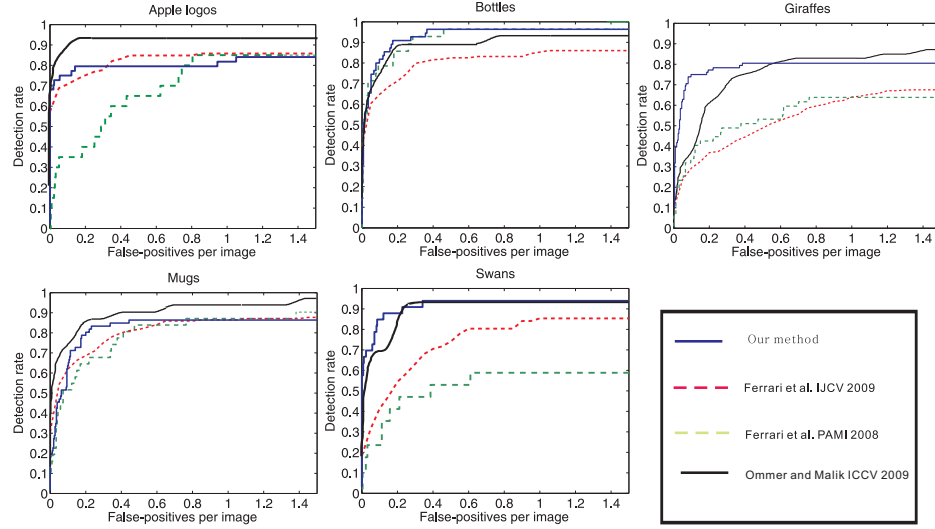


Fig. 5. DR/FPPI curves on ETHZ dataset with comparison to methods in [29, 20, 7].

precision and recall. Similar to [29], both of them do not report actual values and we can only compare by visual estimation. With comparison to [8], it is apparent that our method performs better on swans and equally well on apple logos, mugs and giraffes, but it is outperformed on bottles. Compared to [11], we perform better on giraffes and bottles and we are comparable on classes mugs and swans. The only class we are worse is apple logos.

7.2 Subset of Caltech-101

We also test on two classes of Caltech-101 dataset: cups and car-sides. Since there are no given edge maps like for ETHZ [20], we use Canny Edge detector

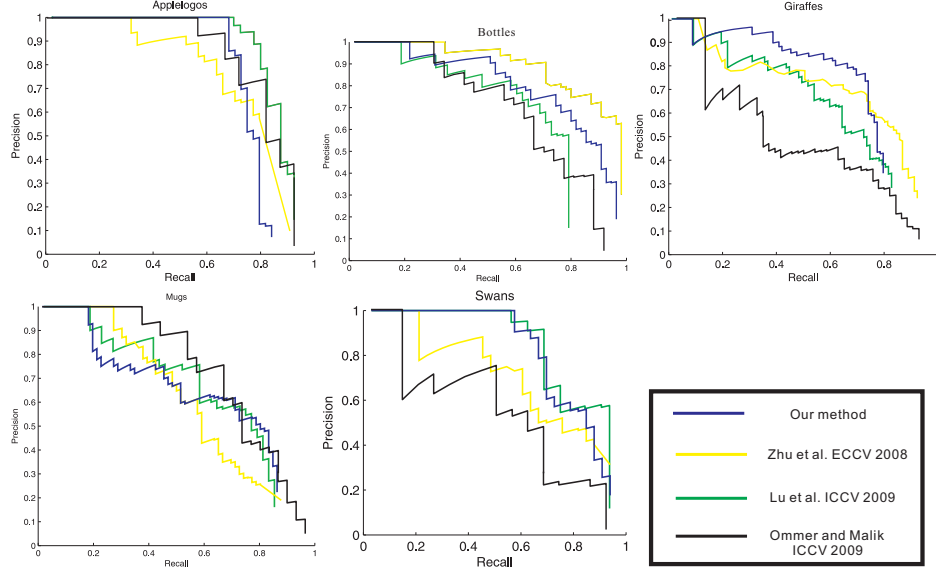


Fig. 6. Precision/Recall curves on ETHZ dataset with comparison to [29, 8, 11].

to obtain edges. Four examples for each class are shown in the last row in Fig. 2.

The blue DR/FPPI curves in Fig 7 illustrate the performance of our algorithm. We obtain detection rate 76.5% and 51.2% at 0.4 FPPI on class cups and car-sides respectively. The high false positive rates for class car-side is due to the low resolution of the images, which makes edge maps not reliable. Consequently, the image segments may be messy making the obtained edge segments different from the model, so that the final global shape similarity may not be able to distinguish the false positives. However, the reasonable detection rate, about 75%, shows that even when the edges are not reliable, the proposed method can still accurately localize the objects.

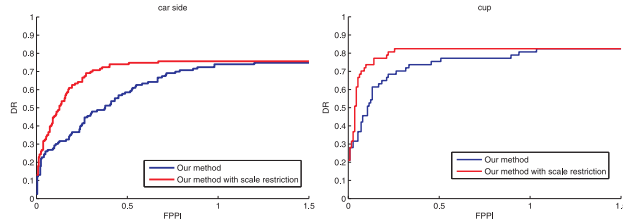


Fig. 7. DR/FPPI curves on class car-side and cups of Caltech-101 dataset.

In order to show the influence of predetermined scale range, we use the scale range as a constraint for our detection results on the two classes. If the scale of a detection hypothesis (which is automatically determined by shape similarity) is not within the scale range, we discard the hypothesis. As shown by red curves in Fig. 7, it is obvious that the restriction of scale range improves the performance. It is mainly due to removing the false positives. With the scale restriction, the detection rate at 0.4 FPPI on class cups and car-sides increases to 82.3% and 73.2% respectively. The increase in detection rate on car-sides is by 22%. We observe that the detection rate of 82.3% on cups is better than the 78.6% reported in [20]. This is particularly impressive, since our contour model for cups is the ETHZ mug model without any modification, while [20] trained a cup classifier on half of cup images.

8 Conclusion

In this paper, we presented a simple yet effective purely shape based approach for object detection. We formulate object detection as a matching problem between image and model segments. To solve the problem of missing segments, we transform the matching problem into finding dominant sets in the correspondence graph. With this transformation, the algorithm is also robust to outlier and noise (background clutter) in the image. Besides, the proposed method can detect multiple objects at multiple scales in one pass, which reduces the complexity a lot compared to standard sliding window and Hough voting approaches.

References

1. Leibe, B., Leonardis, A., Sziele, B.: Combined object categorization and segmentation with an implicit shape model. In: *Proceedings of the Workshop on Statistical Learning in Computer Vision*, Prague, Czech Republic (2004)
2. Gall, J., Lempitsky, V.: Class-specific hough forests for object detection. In: *CVPR*. (2009)
3. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: *CVPR*. (2006)
4. Shotton, J., Winn, J.M., Rother, C., Criminisi, A.: Textonboost: Joint appearance, shape and context modeling for multi-class object recognition and segmentation. In: *ECCV*. (2006)
5. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. In: *CVPR*. (2001)
6. Andreas Opelt, A.P., Zisserman, A.: A boundary-fragment-model for object detection. In: *ECCV*. (2006)
7. Ferrari, V., Jurie, F., Schmid, C.: From images to shape models for object detection. *IJCV* (accepted)
8. Zhu, Q., Wang, L., Wu, Y., Shi, J.: Contour context selection for object detection: a set-to-set contour matching approach. In: *ECCV*. (2008)
9. Stark, M., Goesele, M., Sziele, B.: A shape-based object class model for knowledge transfer. In: *ICCV*. (2009)
10. Jiang, T., Jurie, F., Schmid, C.: Learning shape prior models for object matching. In: *CVPR*. (2009)

11. Lu, C., Latecki, L.J., Adluru, N., Yang, X., Ling, H.: Shape guided contour grouping with particle filters. In: ICCV. (2009)
12. Martin, D., Fowlkes, C., Malik, J.: Learning to detect natural image boundaries using local brightness, color, and texture cues. IEEE PAMI (2004)
13. Rosenfeld, A., Hummel, R., Zucker, S.: Scene labeling by relaxation operations. Trans. on Systems, Man and Cybernetics **6** (1976) 420–433
14. Caetano, T., Caelli, T., Schuurmans, D., Barone, D.: Graphical models and point pattern matching. IEEE PAMI (2006)
15. Berg, A., Berg, T., Malik, J.: Shape matching and object recognition using low distortion correspondences. In: CVPR. (2005)
16. Pavan, M., Pelillo, M.: Dominant sets and pairwise clustering. IEEE. PAMI (2007)
17. Liu, H., Yan, S.: Common visual pattern discovery via spatially coherent correspondences. In: CVPR. (2010)
18. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: CVPR. (2005)
19. Hough, P.: Method and means for recognizing complex patterns. In: U.S. Patent 3069654. (1962)
20. Ferrari, V., Fevrier, L., Jurie, F., Schmid, C.: From images to shape models for object detection. IEEE Trans. PAMI **30** (2008) 36–51
21. Fei-Fei, L., Fergus, R., Perona, P.: One-shot learning of edges and object boundaries. IEEE Trans. on PAMI **28** (2006) 594–611
22. Shotton, J., Blake, A., Cipolla, R.: Multi-scale categorical object recognition using contour fragments. IEEE Trans. on PAMI **30** (2008) 1270–1281
23. Nhon H. Trinh, B.B.K.: Category-specific object recognition and segmentation using a skeletal shape model. In: BMVC. (2009)
24. Bai, X., Wang, X., Latecki, L.J., Tu, Z.: Active skeleton for non-rigid object detection. In: ICCV. (2009)
25. Ravishankar, S., Jain, A., Mittal, A.: Multi-stage contour based detection of deformable objects. In: ECCV. (2008)
26. Belongie, S., Malik, J., Puzicha, J.: Shape matching and object recognition using shape contexts. IEEE Trans. PAMI **24** (2002) 705–522
27. Gu, C., Lim, J.J., Arbelaez, P., Malik, J.: Recognition using regions. In: CVPR. (2009)
28. Maji, S., Malik, J.: Object detection using a max-margin hough transform. In: CVPR. (2009)
29. Ommer, B., Malik, J.: Multi-scale object detection by clustering lines. In: ICCV. (2009)
30. Kovesi, P.D.: Matlab and octave functions for computer vision and image processing. available from <http://www.csse.uwa.edu.au/pk/research/matlabfns/>. (2008)
31. Motzkin, T., Straus, E.: Maxima for graphs and a new proof of a theorem of turan. Canad. J. Math (1965)
32. Sarkar, S., Boyer, K.: Quantitative measures of change based on feature organization: Eigenvalues and eigenvectors. CVIU **71** (1998) 110–136
33. Leordeanu, M., Hebert, M.: A spectral technique for correspondence problems using pairwise constraints. In: ICCV. (2005)
34. Weibull, J.: Evolutionary game theory. MIT Press (1997)
35. Everingham, M., Gool, L.V., Williams, C.K.I., Winn, J., Zisserman, A.: (The pascal visual object classes challenge 2008 (voc 2008) results. <http://pascallin.ecs.soton.ac.uk/challenges/voc/voc2008/workshop/>)