

# Multiscale Combinatorial Grouping for Image Segmentation and Object Proposal Generation

Jordi Pont-Tuset\*, Pablo Arbeláez\*, Jonathan T. Barron, *Member, IEEE*,  
Ferran Marques, *Senior Member, IEEE*, Jitendra Malik, *Fellow, IEEE*

**Abstract**—We propose a unified approach for bottom-up hierarchical image segmentation and object proposal generation for recognition, called Multiscale Combinatorial Grouping (MCG). For this purpose, we first develop a fast normalized cuts algorithm. We then propose a high-performance hierarchical segmenter that makes effective use of multiscale information. Finally, we propose a grouping strategy that combines our multiscale regions into highly-accurate object proposals by exploring efficiently their combinatorial space. We also present Single-scale Combinatorial Grouping (SCG), a faster version of MCG that produces competitive proposals in under five second per image. We conduct an extensive and comprehensive empirical validation on the BSDS500, SegVOC12, SBD, and COCO datasets, showing that MCG produces state-of-the-art contours, hierarchical regions, and object proposals.

**Index Terms**—Image segmentation, object proposals, normalized cuts.

## 1 INTRODUCTION

TWO paradigms have shaped the field of object recognition in the last decade. The first one, popularized by the Viola-Jones face detection algorithm [1], formulates object localization as window classification. The basic scanning-window architecture, relying on histograms of gradients and linear support vector machines, was introduced by Dalal and Triggs [2] in the context of pedestrian detection and is still at the core of seminal object detectors on the PASCAL challenge such as Deformable Part Models [3].

The second paradigm relies on perceptual grouping to provide a limited number of high-quality and category-independent object proposals, which can then be described with richer representations and used as input to more sophisticated learning methods. Examples in this family are [4], [5]. Recently, this approach has dominated the PASCAL segmentation challenge [6], [7], [8], [9], improved object detection [10], fine-grained categorization [11] and proven competitive in large-scale classification [12].

Since the power of this second paradigm is critically dependent on the accuracy and the number of object proposals, an increasing body of research has delved



Fig. 1. **Top**: original image, instance-level ground truth from COCO and our multiscale hierarchical segmentation. **Bottom**: our best object proposals among 150.

into the problem of their generation [13], [14], [15], [12], [16], [17], [18], [19]. However, those approaches typically focus on learning generic properties of objects from a set of examples, while reasoning on a fixed set of regions and contours produced by external bottom-up segmenters such as [20], [21].

In this paper, we propose a unified approach to multiscale hierarchical segmentation and object proposal generation called Multiscale Combinatorial Grouping (MCG). Fig. 1 shows an example of our results and Fig. 2 an overview of our pipeline. Our main contributions are:

- An efficient normalized cuts algorithm, which in practice provides a  $20\times$  speed-up to the eigenvector computation required for contour globalization [20], [22] (Sect. 3.1).
- A state-of-the-art hierarchical segmenter that leverages multiscale information (Sect. 3.3).
- A grouping algorithm that produces accurate object proposals by efficiently exploring the combinatorial space of our multiscale regions (Sect. 5).

• J. Pont-Tuset and F. Marques are with the Department of Signal Theory and Communications, Universitat Politècnica de Catalunya, BarcelonaTech (UPC), Spain. E-mail: {jordi.pont,ferran.marques}@upc.edu  
• P. Arbeláez is with the Department of Biomedical Engineering, Universidad de los Andes, Colombia. E-mail: pa.arbelaez@uniandes.edu.co  
• J. T. Barron, and J. Malik are with the Department of Electrical Engineering and Computer Science, University of California at Berkeley, Berkeley, CA 94720. E-mail: {barron,malik}@eecs.berkeley.edu

\* The first two authors contributed equally

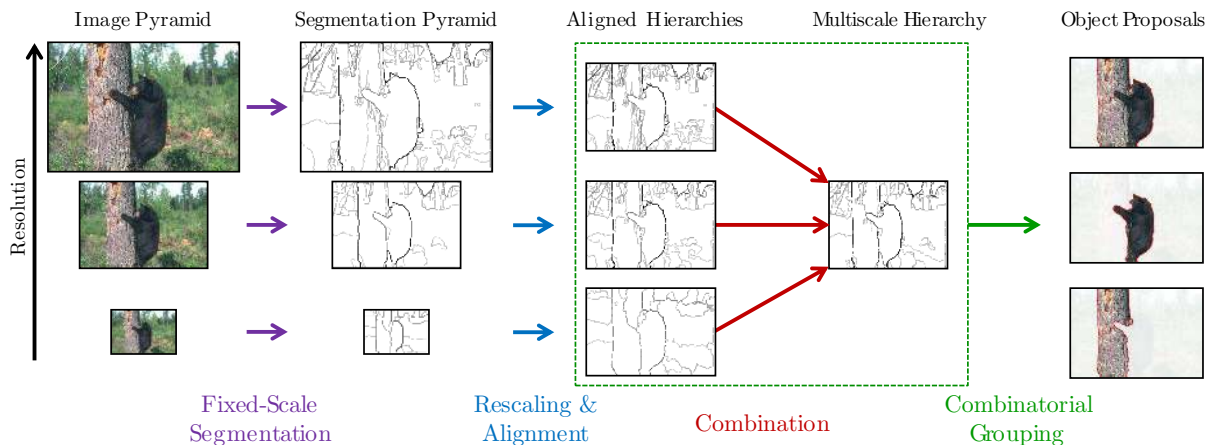


Fig. 2. **Multiscale Combinatorial Grouping.** Starting from a multiresolution image pyramid, we perform hierarchical segmentation at each scale independently. We align these multiple hierarchies and combine them into a single multiscale segmentation hierarchy. Our grouping component then produces a ranked list of object proposals by efficiently exploring the combinatorial space of these regions.

We conduct a comprehensive and large-scale empirical validation. On the BSDS500 (Sect. 4) we report significant progress in contour detection and hierarchical segmentation. On the VOC2012, SBD, and COCO segmentation datasets (Sect. 6), our proposals obtain overall state-of-the-art accuracy both as segmented proposals and as bounding boxes. MCG is efficient, its good generalization power makes it parameter free in practice, and it provides a ranked set of proposals that are competitive in all regimes of number of proposals.

## 2 RELATED WORK

For space reasons, we focus our review on recent normalized cut algorithms and object proposals for recognition.

**Fast normalized cuts:** The efficient computation of normalized-cuts eigenvectors has been the subject of recent work, as it is often the computational bottleneck in grouping algorithms. Taylor [23] presented a technique for using a simple watershed oversegmentation to reduce the size of the eigenvector problem, sacrificing accuracy for speed. We take a similar approach of solving the eigenvector problem in a reduced space, though we use simple image-pyramid operations on the affinity matrix (instead of a separate segmentation algorithm) and we see no loss in performance despite a  $20\times$  speed improvement. Maire and Yu [24] presented a novel multigrid solver for producing eigenvectors at multiple scales, which speeds up fine-scale eigenvector computation by leveraging coarse-scale solutions. Our technique also uses the scale-space structure of an image, but instead of solving the problem at multiple scales, we simply reduce the scale of the problem, solve it at a reduced scale, and then upsample the solution while preserving the structure of the image. As such, our technique is faster and much simpler, requiring only a few lines of code wrapped around a standard sparse eigensolver.

**Object Proposals:** Class-independent methods that generate object hypotheses can be divided into those whose output is an image window and those that generate segmented proposals.

Among the former, Alexe *et al.* [16] propose an *objectness* measure to score randomly-sampled image windows based on low-level features computed on the superpixels of [21]. Manen *et al.* [25] propose to use the Randomized Prim’s algorithm, Zitnick *et al.* [26] group contours directly to produce object windows, and Cheng *et al.* [27] generate box proposals at 300 images per second. In contrast to these approaches, we focus on the finer-grained task of pixel-accurate object extraction, rather than on window selection. However, by just taking the bounding box around our segmented proposals, our results are also state of the art as window proposals.

Among the methods that produce segmented proposals, Carreira and Sminchisescu [18] hypothesize a set of placements of fore- and background seeds and, for each configuration, solve a constrained parametric min-cut (CPMC) problem to generate a pool of object hypotheses. Endres and Hoiem [19] base their category-independent object proposals on an iterative generation of a hierarchy of regions, based on the contour detector of [20] and occlusion boundaries of [28]. Kim and Grauman [17] propose to match parts of the shape of exemplar objects, regardless of their class, to detected contours by [20]. They infer the presence and shape of a proposal object by adapting the matched object to the computed superpixels.

Uijlings *et al.* [12] present a selective search algorithm based on segmentation. Starting with the superpixels of [21] for a variety of color spaces, they produce a set of segmentation hierarchies by region merging, which are used to produce a set of object proposals. While we also take advantage of different hierarchies to gain diversity, we leverage multiscale information rather than different

color spaces.

Recently, two works proposed to train a cascade of classifiers to learn which sets of regions should be merged to form objects. Ren and Shankhnavich [29] produce full region hierarchies by iteratively merging pairs of regions and adapting the classifiers to different scales. Weiss and Taskar [30] specialize the classifiers also to size and class of the annotated instances to produce object proposals.

Malisiewicz and Efros [4] took one of the first steps towards combinatorial grouping, by running multiple segmenters with different parameters and merging up to three adjacent regions. In [8], another step was taken by considering hierarchical segmentations at three different scales and combining pairs and triplets of adjacent regions from the two coarser scales to produce object proposals.

The most recent wave of object proposal algorithms is represented by [13], [14], and [15], which all keep the quality of the seminal proposal works while improving the speed considerably. Krähenbühl and Koltun [13] find object proposal by identifying critical level sets in geodesic distance transforms, based on seeds placed in learnt places in the image. Rantalankila *et al.* [14] perform a global and local search in the space of sets of superpixels. Humayun *et al.* [15] reuse a graph to perform many parametric min-cuts over different seeds in order to speed the process up.

A substantial difference between our approach and previous work is that, instead of relying on pre-computed hierarchies or superpixels, we propose a unified approach that produces and groups high-quality multiscale regions. With respect to the combinatorial approaches of [4], [8], our main contribution is to develop efficient algorithms to explore a much larger combinatorial space by taking into account a set of object examples, increasing thus the likelihood of having complete objects in the pool of proposals. Our approach has therefore the flexibility to adapt to specific applications and types of objects, and can produce proposals at any trade-off between their number and their accuracy.

### 3 THE SEGMENTATION ALGORITHM

Consider a segmentation of the image into regions that partition its domain  $\mathcal{S} = \{S_i\}_i$ . A segmentation hierarchy is a family of partitions  $\{S^*, S^1, \dots, S^L\}$  such that: (1)  $S^*$  is the finest set of *superpixels*, (2)  $S^L$  is the complete domain, and (3) regions from coarse levels are unions of regions from fine levels. A hierarchy where each level  $S^i$  is assigned a real-valued index  $\lambda_i$  can be represented by a dendrogram, a region tree where the height of each node is its index. Furthermore, it can also be represented as an ultrametric contour map (UCM), an image obtained by weighting the boundary of each pair of adjacent regions in the hierarchy by the index at which they are merged [31], [32]. This representation unifies the problems of contour detection and hierarchical image segmentation:

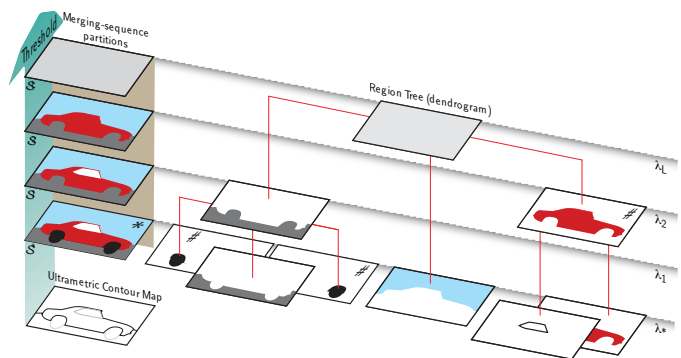


Fig. 3. **Duality between a UCM and a region tree:** Schematic view of the dual representation of a segmentation hierarchy as a region dendrogram and as an ultrametric contour map.

a threshold at level  $\lambda_i$  in the UCM produces the segmentation  $S^i$ .

Figure 3 schematizes these concepts. First, the lower left corner shows the probability of boundary of a UCM. One of the main properties of a UCM is that when we threshold the contour strength at a certain value, we obtain a closed boundary map, and thus a partition. Thresholding at different  $\lambda_i$ , therefore, we obtain the so-called merging-sequence partitions (left column in Figure 3); named after the fact that a step in this sequence corresponds to merging the set of regions sharing the boundary of strength exactly  $\lambda_i$ .

For instance, the boundary between the wheels and the floor has strength  $\lambda_1$ , thus thresholding the contour above  $\lambda_1$  makes the wheels *merge* with the floor. If we represent the regions in a partition as nodes of a graph, we can then represent the result of merging them as their parent in a tree. The result of sweeping all  $\lambda_i$  values can therefore be represented as a region tree, whose root is the region representing the whole image (right part of Figure 3). Given that each *merging* is associated with a contour strength, the region tree is in fact a region dendrogram.

As an example, in the gPb-ucm algorithm of [20], brightness, color and texture gradients at three fixed disk sizes are first computed. These local contour cues are globalized using spectral graph-partitioning, resulting in the gPb contour detector. Hierarchical segmentation is then performed by iteratively merging adjacent regions based on the average gPb strength on their common boundary. This algorithm produces therefore a tree of regions at multiple levels of homogeneity in brightness, color and texture, and the boundary strength of its UCM can be interpreted as a measure of contrast.

Coarse-to-fine is a powerful processing strategy in computer vision. We exploit it in two different ways to develop an efficient, scalable and high-performance segmentation algorithm: (1) To speed-up spectral graph partitioning and (2) To create aligned segmentation hierarchies.

### 3.1 Fast Downsampled Eigenvector Computation

The normalized cuts criterion is a key globalization mechanism of recent high-performance contour detectors such as [20], [22]. Although powerful, such spectral graph partitioning has a significant computational cost and memory footprint that limit its scalability. In this section, we present an efficient normalized cuts approximation which in practice preserves full performance for contour detection, while having low memory requirements and providing a  $20\times$  speed-up.

Given a symmetric affinity matrix  $A$ , we would like to compute the  $k$  smallest eigenvectors of the Laplacian of  $A$ . Directly computing such eigenvectors can be very costly even with sophisticated solvers, due to the large size of  $A$ . We therefore present a technique for efficiently approximating the eigenvector computation by taking advantage of the multiscale nature of our problem:  $A$  models affinities between pixels in an image, and images naturally lend themselves to multiscale or pyramid-like representations and algorithms.

Our algorithm is inspired by two observations: 1) if  $A$  is bistochastic (the rows and columns of  $A$  sum to 1) then the eigenvectors of the Laplacian  $A$  are equal to the eigenvectors of the Laplacian of  $A^2$ , and 2) because of the scale-similar nature of images, the eigenvectors of a “downsampled” version of  $A$  in which every other pixel has been removed should be similar to the eigenvectors of  $A$ . Let us define `pixel_decimate( $A$ )`, which takes an affinity matrix  $A$  and returns the set of indices of rows/columns in  $A$  corresponding to a decimated version of the image from which  $A$  was constructed. That is, if  $i = \text{pixel\_decimate}(A)$ , then  $A[i, i]$  is a decimated matrix in which alternating rows and columns of the image have been removed. Computing the eigenvectors of  $A[i, i]$  works poorly, as decimation disconnects pixels in the affinity matrix, but the eigenvectors of the decimated squared affinity matrix  $A^2[i, i]$  are similar to those of  $A$ , because by squaring the matrix before decimation we intuitively allow each pixel to propagate information to all of its neighbors in the graph, maintaining connections even after decimation. Our algorithm works by efficiently computing  $A^2[i, i]$  as  $A[:, i]^T A[:, i]$  (the naive approach of first squaring  $A$  and then decimating it is prohibitively expensive), computing the eigenvectors of  $A^2[i, i]$ , and then “upsampling” those eigenvectors back to the space of the original image by pre-multiplying by  $A[:, i]$ . This squaring-and-decimation procedure can be applied recursively several times, each application improving efficiency while slightly sacrificing accuracy.

Pseudocode for our algorithm, which we call “DNCuts” (Downsampled Normalized Cuts) is given in Algorithm 1, where  $A$  is our affinity matrix and  $d$  is the number of times that our squaring-and-decimation operation is applied. Our algorithm repeatedly applies our joint squaring-and-decimation procedure, computes

---

#### Algorithm 1 `dncuts( $A, d, k$ )`

---

```

1:  $A_0 \leftarrow A$ 
2: for  $s = [1, 2, \dots, d]$  do
3:    $i_s \leftarrow \text{pixel\_decimate}(A_{s-1})$ 
4:    $B_s \leftarrow A_{s-1}[:, i_s]$ 
5:    $C_s \leftarrow \text{diag}(B_s \mathbf{1})^{-1} B_s$ 
6:    $A_s \leftarrow C_s^T B_s$ 
7: end for
8:  $X_d \leftarrow \text{ncuts}(A_d, k)$ 
9: for  $s = [d, d-1, \dots, 1]$  do
10:   $X_{s-1} \leftarrow C_s X_s$ 
11: end for
12: return whiten( $X_0$ )

```

---

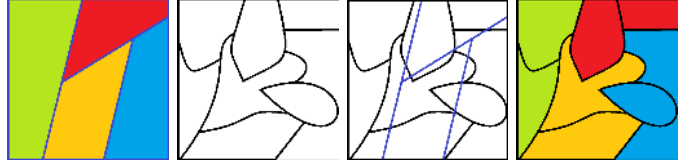


Fig. 4. **Example of segmentation projection.** In order to “snap” the boundaries of a segmentation  $\mathcal{R}$  (left) to those of a segmentation  $\mathcal{S}$  (middle), since they do not align, we compute  $\pi(\mathcal{R}, \mathcal{S})$  (right) by assigning to each segment in  $\mathcal{S}$  its mode among the labels of  $\mathcal{R}$ .

the smallest  $k$  eigenvectors of the final “downsampled” matrix  $A_d$  by using a standard sparse eigensolver `ncuts( $A_d, k$ )`, and repeatedly “upsamples” those eigenvectors. Because our  $A$  is not bistochastic and decimation is not an orthonormal operation, we must do some normalization throughout the algorithm (line 5) and whiten the resulting eigenvectors (line 10). We found that values of  $d = 2$  or  $d = 3$  worked well in practice. Larger values of  $d$  yielded little speed improvement (as much of the cost is spent downsampling  $A_0$ ) and start negatively affecting accuracy. Our technique is similar to Nystrom’s method for computing the eigenvectors of a subset of  $A$ , but our squaring-and-decimation procedure means that we do not depend on long-range connections between pixels.

### 3.2 Aligning Segmentation Hierarchies

In order to leverage multi-scale information, our approach combines segmentation hierarchies computed independently at multiple image resolutions. However, since subsampling an image removes details and smooths away boundaries, the resulting UCMs are misaligned, as illustrated in the second panel of Fig. 2. In this section, we propose an algorithm to align an arbitrary segmentation hierarchy to a target segmentation and, in Sect. 5, we show its effectiveness for multi-scale segmentation.

The basic operation is to “snap” the boundaries of a segmentation  $\mathcal{R} = \{R_i\}_i$  to a segmentation  $\mathcal{S} = \{S_j\}_j$ , as illustrated in Fig. 4. For this purpose, we define  $\mathcal{L}(S_j)$ ,

1. The *Matlab-like* notation  $A[:, i]$  indicates keeping the columns of matrix  $A$  whose indices are in the set  $i$ .

the new label of a region  $S_j \in \mathcal{S}$ , as the majority label of its pixels in  $\mathcal{R}$ :

$$\mathcal{L}(S_j) = \arg \max_i \frac{|S_j \cap R_i|}{|S_j|} \quad (1)$$

We call the segmentation defined by this new labeling of all the regions of  $\mathcal{S}$  the *projection* of  $\mathcal{R}$  onto  $\mathcal{S}$  and denote it by  $\pi(\mathcal{R}, \mathcal{S})$ .

In order to project an UCM onto a target segmentation  $\mathcal{S}$ , which we denote  $\pi(\text{UCM}, \mathcal{S})$ , we project in turn each of the levels of the hierarchy onto  $\mathcal{S}$ . Note that, since all the levels are projected onto the same segmentation, the family of projections is by construction a hierarchy of segmentations. This procedure is summarized in pseudo-code in Algorithm 2.

---

#### Algorithm 2 UCM Rescaling and Alignment

---

**Require:** An UCM with a set of levels  $[t_1, \dots, t_K]$

**Require:** A target segmentation  $\mathcal{S}^*$

```

1:  $\text{UCM}_\pi \leftarrow 0$ 
2: for  $t = [t_1, \dots, t_K]$  do
3:    $\mathcal{S} \leftarrow \text{sampleHierarchy}(\text{UCM}, t)$ 
4:    $\mathcal{S} \leftarrow \text{rescaleSegmentation}(\mathcal{S}, \mathcal{S}^*)$ 
5:    $\mathcal{S} \leftarrow \pi(\mathcal{S}, \mathcal{S}^*)$ 
6:    $\text{contours} \leftarrow \text{extractBoundary}(\mathcal{S})$ 
7:    $\text{UCM}_\pi \leftarrow \max(\text{UCM}_\pi, t * \text{contours})$ 
8: end for
9: return  $\text{UCM}_\pi$ 

```

---

Observe that the routines `sampleHierarchy` and `extractBoundary` can be computed efficiently because they involve only thresholding operations and connected components labeling. The complexity is thus dominated by `rescaleSegmentation` in Step 4, a nearest neighbor interpolation, and the projection in Step 5, which are computed  $K$  times.

### 3.3 Multiscale Hierarchical Segmentation

**Single-scale segmentation:** We consider as input the following local contour cues: (1) brightness, color and texture differences in half-disks of three sizes [33], (2) sparse coding on patches [22], and (3) structured forest contours [34]. We globalize the contour cues independently using our fast eigenvector gradients of Sect. 3.1, combine global and local cues linearly, and construct an UCM based on the mean contour strength. We tried learning weights using gradient ascent on the F-measure on the training set [20], but evaluating the final hierarchies rather than open contours. We observed that this objective favors the quality of contours at the expense of regions and obtained better overall results by optimizing the Segmentation Covering metric [20].

**Hierarchy Alignment:** We construct a multiresolution pyramid with  $N$  scales by subsampling / super-sampling the original image and applying our single-scale segmenter. In order to preserve thin structures and details, we declare as set of possible boundary

locations the finest superpixels in the highest-resolution. Then, applying recursively Algorithm 2, we project each coarser UCM onto the next finer scale until aligning it to the highest resolution superpixels.

**Multiscale Hierarchy:** After alignment, we have a fixed set of boundary locations, and  $N$  strengths for each of them, coming from the different scales. We formulate this problem as binary boundary classification and train a classifier that combines these  $N$  features into a single probability of boundary estimation. We experimented with several learning strategies for combining UCM strengths: (a) Uniform weights transformed into probabilities with Platt’s method. (b) SVMs and logistic regression, with both linear and additive kernels. (c) Random Forests. (d) The same algorithm as for single-scale. We found the results with all learning methods surprisingly similar, in agreement with the observation reported by [33]. This particular learning problem, with only a handful of dimensions and millions of data points, is relatively easy and performance is mainly driven by our already high-performing and well calibrated features. We therefore use the simplest option (a).

## 4 EXPERIMENTS ON THE BSDS500

We conduct extensive experiments on the BSDS500 [35], using the standard evaluation metrics and following the best practice rules of that dataset. We also report results with a recent evaluation metric  $F_{op}$  [36], [37], Precision-Recall for objects and parts, using the publicly-available code.

**Single-scale Segmentation:** Table 1-top shows the performance of our single-scale segmenter for different types of input contours on the validation set of the BSDS500. We obtain high-quality hierarchies for all the cues considered, showing the generality of our approach. Furthermore, when using them jointly (row ‘Comb.’ in top panel), our segmenter outperforms the versions with individual cues, suggesting its ability to leverage diversified inputs. In terms of efficiency, our fast normalized cuts algorithm provides an average 20× speed-up over [20], starting from the same local cues, with no significant loss in accuracy and with a low memory footprint.

**Multiscale Segmentation:** Table 1-bottom evaluates our full approach in the same experimental conditions as the upper panel. We observe a consistent improvement in performance in all the metrics for all the inputs, which validates our architecture for multiscale segmentation. We experimented with the range of scales and found  $N = \{0.5, 1, 2\}$  adequate for our purposes. A finer sampling or a wider range of scales did not provide noticeable improvements. We tested also two degraded versions of our system (not shown in the table). For the first one, we resized contours to the original image resolution, created UCMs and combined them with the same method as our final system. For the second one, we transformed per-scale UCMs to the original resolution,

	Input	Boundary		Region							
		$F_b$		$F_{op}$		SC		PRI		VI	
		ODS	OIS	ODS	OIS	ODS	OIS	ODS	OIS	ODS	OIS
Single-Scale	Pb [33]	0.702	0.733	0.334	0.370	0.577	0.636	0.801	0.847	1.692	1.490
	SC [22]	0.697	0.725	0.264	0.306	0.540	0.607	0.777	0.835	1.824	1.659
	SF [34]	0.719	0.737	0.338	0.399	0.582	0.651	0.803	0.851	1.608	1.432
	Comb.	0.719	0.750	0.358	0.403	0.602	0.655	0.809	0.855	1.619	1.405
Multiscale	Pb [33]	0.713	0.745	0.350	0.389	0.598	0.656	0.807	0.856	1.601	1.418
	SC [22]	0.705	0.734	0.331	0.384	0.579	0.647	0.799	0.851	1.637	1.460
	SF [34]	0.725	0.744	0.370	0.420	0.600	0.660	0.810	0.854	1.557	1.390
	Comb.	0.725	0.757	0.371	0.408	0.611	0.670	0.813	0.862	1.548	1.367

TABLE 1

**BSDS500 val set.** Control experiments for single-scale (top) and multiscale (bottom) hierarchical segmentation with different input contour detectors

but omitted the strength transfer to the finest superpixels before combining them. The first ablated version produces interpolation artifacts and smooths away details, while the second one suffers from misalignment. Both fail to improve performance over the single-scale result, which provides additional empirical support for our multiscale approach. We also observed a small degradation in performance when forcing the input contour detector to use only the original image resolution, which indicates the advantages of considering multiscale information at all stages of processing.

Since there are no drastic changes in our results when taking as input the different individual cues or their combination, in the sequel we use the version with structured forests for efficiency reasons, which we denote *MCG-UCM-Our*.

**Comparison with state-of-the-art.:** Figure 5 compares our multiscale hierarchical segmenter MCG (—●—) and our single-scale hierarchical segmenter SCG (—○—) on the BSDS500 test set against all the methods for which there is publicly available code. We also compare to the recent ISCRA [29] hierarchies (—→—), provided precomputed by the authors. We obtain consistently the best results to date on BSDS500 for all operating regimes, both in terms of boundary and region quality.

Note that the range of object scales in the BSDS500 is limited, which translates into modest absolute gains from MCG (—●—) with respect to SCG (—○—) in terms of boundary evaluation (left-hand plot), but more significant improvements in terms of objects and parts (right-hand plot). We will also observe more substantial improvements with respect to gPb-UCM (—▲—) when we move to PASCAL, SBD, and COCO in Section 6 (e.g. see Fig. 9).

**Ground-Truth Hierarchy:** In order to gain further insights, we transfer the strength of each ground-truth segmentation to our highest-resolution superpixels  $S^{N^*}$  and construct a combined hierarchy. This approximation to the semantic hierarchy, Ground-Truth Hierarchy (GTH) in Fig. 5, is an upper-bound for our approach as both share the same boundary locations and the only difference is their strength. Since the strength of GTH

is proportional to the number of subjects marking it, it provides naturally the correct semantic ordering, where outer object boundaries are stronger than internal parts.

Recently, Maire *et al.* [38] developed an annotation tool where the user encodes explicitly the “perceptual strength” of each contour. Our approach provides an alternative where the semantic hierarchy is reconstructed by sampling flat annotations from multiple subjects.

## 5 OBJECT PROPOSAL GENERATION

The image segmentation algorithm presented in the previous sections builds on low-level features, so its regions are unlikely to represent accurately complete objects with heterogeneous parts. In this context, object proposal techniques create a set of hypotheses, possibly overlapping, which are more likely to represent full object instances.

Our approach to object proposal generation is to combinatorially look for sets of regions from our segmentation hierarchies that merged together are likely to represent complete objects. In this section, we first describe the efficient computation of certain region descriptors on a segmentation tree. Then, we describe how we use these techniques to efficiently explore the sets of merged regions from the hierarchy. Finally, we explain how we train the parameters of our algorithm for object proposals and how we rank the candidates by their probability of representing an object.

**Fast computation of descriptors:** Let us assume, for instance, we want to compute the area of all regions in the hierarchy. Intuitively, working strictly on the merging-sequence partitions, we would need to scan all pixels in all partitions. On the other hand, working on the region tree allows us to scan the image only once to compute the area of the leaves, and then propagate the area to all parents as the addition of the areas of their children.

As a drawback, the algorithms become intricate in terms of coding and necessary data structures. Take, for instance, the computation of the neighbors of a certain region, which is trivial via scanning the partition on the merging sequence (look for region labels in the adjacent boundary pixels), but need tailored data structures and algorithms in the region tree.

Formally, let us assume the image has  $p$  pixels, and we build a hierarchy based on  $s$  superpixels (leaves of the tree), and  $m$  mergings (different UCM strength values). The cost of computing the area on all regions using the merging-sequence partitions will be the cost of scanning all pixels in these partitions, thus  $p \cdot (m+1)$ . In contrast, the cost using the region tree will involve scanning the image once, and then propagating the area values, so  $p+m$ , which is notably faster.

We built tailored algorithms and data structures to compute the bounding box, perimeter, and neighbors of a region using the region tree representation.

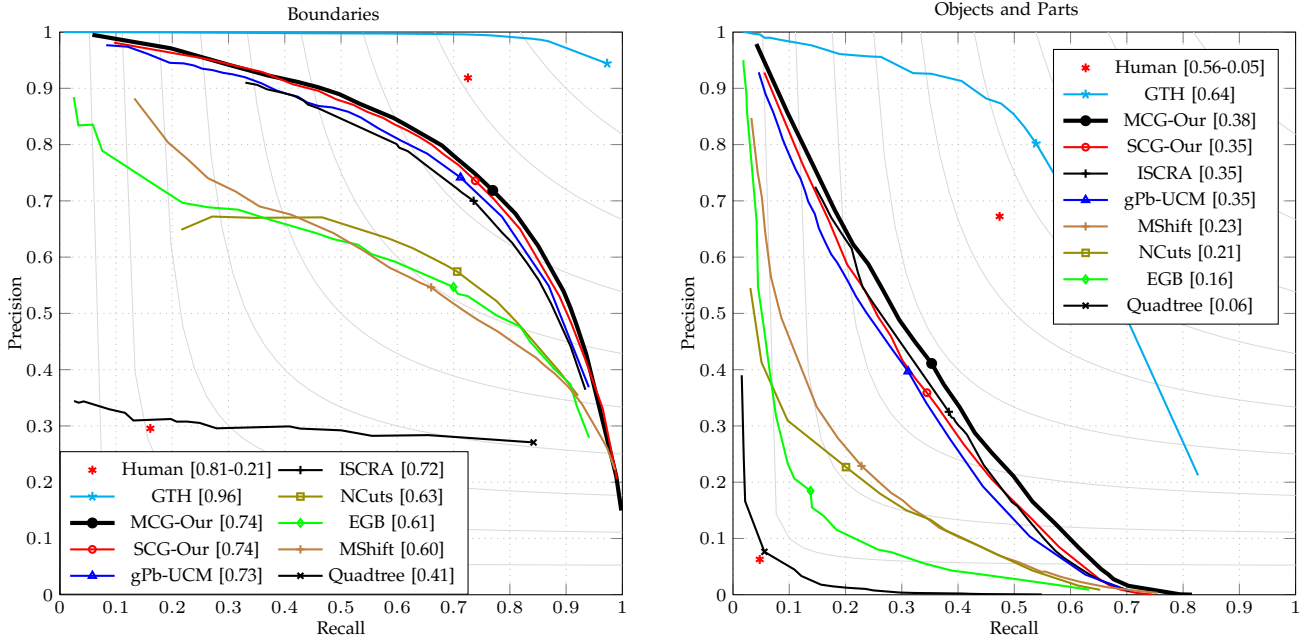


Fig. 5. **BSDS500 test set.** Precision-Recall curves for boundaries [35] (left) and for objects and parts [36] (right). The marker on each curve is placed on the Optimal Dataset Scale (ODS), and its F measure is presented in brackets in the legend. The isolated red asterisks refer to the human performance assessed on the same image (one human partition against the rest of human annotations on the same image) and on a different image (one human partition against the human partitions of a different, randomly selected, image).

**Combinatorial Grouping of Proposals:** We can cast object segmentation as *selecting* some regions in the hierarchy, or in other words, as a combinatorial optimization problem on the hierarchy. To illustrate this principle, Figure 6(a) shows the simplified representation of the hierarchy in Figure 3. Figure 6(b) and (c) show two object proposals, and their representation in the region hierarchy.

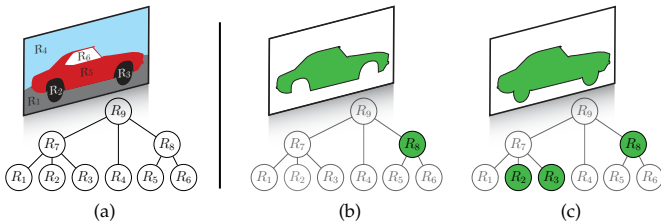


Fig. 6. **Object segmentation as combinatorial optimization:** Examples of objects (b), (c), formed by selecting regions from a hierarchy (a).

Since hierarchies are built taking only low-level features into account, and do not use semantic information, objects will usually not be optimally represented using a single region in the hierarchy. As an example, Figure 6(c) shows the optimum representation of the car, consisting of three regions.

A sensible approach to create object proposals is therefore to explore the set of  $n$ -tuples of regions. The main idea behind MCG is to explore this set efficiently, taking advantage of the region tree representation, via the fast

computation of region neighbors.

The whole set of tuples, however, is huge, and so it is not feasible to explore it exhaustively. Our approach ranks the proposals using the height of their regions in the tree (UCM strength) and explores the tree from the top, but only down to a certain threshold. To speed the process up, we design a top-down algorithm to compute the region neighbors and thus only compute them down to a certain depth in the tree.

To further improve the results, we not only consider the  $n$ -tuples from the resulting MCG-UCM-Our hierarchy, but also the rest of hierarchies computed at different scales. As we will show in the experiments, diversity significantly improves the proposal results.

**Parameter Learning via Pareto Front Optimization:** MCG takes a set of diverse hierarchies and computes the  $n$ -tuples up to a certain UCM strength. We can interpret the  $n$ -tuples from each hierarchy as a ranked list of  $N_i$  proposals that are put together to create the final set of  $N_p$  proposals.

At training time, we would like to find, for different values of  $N_p$ , the number of proposals from each ranked list  $\bar{N}_i$  such that the joint pool of  $N_p$  proposals has the best achievable quality. We frame this learning problem as a Pareto front optimization [39], [40] with two conflicting objective functions: number of proposals and achievable quality. At test time, we select a working point on the Pareto front, represented by the  $\{\bar{N}_i\}$  values, based either on the number of proposals  $N_p$  we can handle or on the minimum achievable quality our

application needs, and we combine the  $\bar{N}_i$  top proposals from each hierarchy list.

Formally, assuming  $R$  ranked lists  $L_i$ , an exhaustive learning algorithm would consider all possible values of the  $R$ -tuple  $\{N_1, \dots, N_R\}$ , where  $N_i \in \{0, \dots, |L_i|\}$ ; adding up to  $\prod_1^R |L_i|$  parameterizations to try, which is intractable in our setting.

Figure 7 illustrates the learning process. To reduce the dimensionality of the search space, we start by selecting two ranked lists  $L_1, L_2$  (green curves) and we sample the list at  $S$  levels of number of proposals (green dots). We then scan the full  $S^2$  different parameterizations to combine the proposals from both (blue dots). In other words, we analyze the sets of proposals created by combining the top  $N_1$  from  $L_1$  (green dots) and the top  $N_2$  from  $L_2$ .

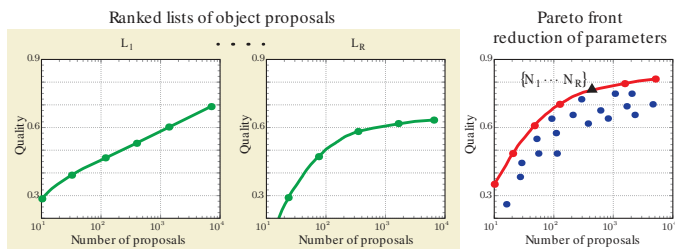


Fig. 7. **Pareto front learning:** Training the combinatorial generation of proposals using the Pareto front

The key step of the optimization consists in discarding those parameterizations whose quality point is not in the Pareto front (red curve). (i.e., those parameterizations that can be substituted by another with better quality with the same number of proposals, or by one with the same quality with less proposals.) We sample the Pareto front to  $S$  points and we iterate the process until all the ranked lists are combined.

Each point in the final Pareto front corresponds to a particular parameterization  $\{N_1, \dots, N_R\}$ . At train time, we choose a point on this curve, either at a given number of proposals  $N_c$  or at the achievable quality we are interested in (black triangle) and store the parameters  $\{\bar{N}_1, \dots, \bar{N}_R\}$ . At test time, we combine the  $\{\bar{N}_1, \dots, \bar{N}_R\}$  top proposals from each ranked list. The number of sampled configurations using the proposed algorithm is  $(R - 1)S^2$ , that is, we have reduced an exponential problem ( $S^R$ ) to a quadratic one.

**Regressed Ranking of Proposals:** To further reduce the number of proposals, we train a regressor from low-level features, as in [18]. Since the proposals are all formed by a set of regions from a reduced set of hierarchies, we focus on features that can be computed efficiently in a bottom-up fashion, as explained previously.

We compute the following features:

- **Size and location:** Area and perimeter of the candidate; area, position, and aspect ratio of the bounding box; and the area balance between the regions in the

candidate.

- **Shape:** Perimeter (and sum of contour strength) divided by the squared root of the area; and area of the region divided by that of the bounding box.

- **Contours:** Sum of contour strength at the boundaries, mean contour strength at the boundaries; minimum and maximum UCM threshold of appearance and disappearance of the regions forming the candidate.

We train a Random Forest using these features to regress the object overlap with the ground truth, and diversify the ranking based on Maximum Marginal Relevance measures [18]. We tune the random forest learning on half training set and validating on the other half. For the final results, we train on the training set and evaluate our proposals on the validation set of PASCAL 2012.

## 6 EXPERIMENTS ON PASCAL VOC, SBD, AND COCO

This section presents our large-scale empirical validation of the object proposal algorithm described in the previous section. We perform experiments in three annotated databases, with a variety of measures that demonstrate the state-of-the-art performance of our algorithm.

**Datasets and Evaluation Measures:** We conduct experiments in the following three annotated datasets: the segmentation challenge of PASCAL 2012 Visual Object Classes (SegVOC12) [41], the Berkeley Semantic Boundaries Dataset (SBD) [42], and the Microsoft Common Objects in Context (COCO) [43]. They all consist of images with annotated objects of different categories. Table 2 summarizes the number of images and object instances in each database.

	Number of Classes	Number of Images	Number of Objects
SegVOC12	20	2 913	9 847
SBD	20	12 031	32 172
COCO	80	123 287	910 983

TABLE 2  
Sizes of the databases

Regarding the performance metrics, we measure the achievable quality with respect to the number of proposals, that is, the quality we would have if an oracle selected the best proposal among the pool. This aligns with the fact that object proposals are a preprocessing step for other algorithms that will represent and classify them. We want, therefore, the achievable quality within the proposals to be as high as possible, while reducing the number of proposals to make the final system as fast as possible.

As a measure of quality of a specific proposal with respect to an annotated object, we consider the Jaccard index  $J$ , also known as overlap or intersection over union; which is defined as the size of the intersection of the two pixel sets over the size of their union.



To compute the overall quality for the whole database, we first select the best proposal for each annotated instance with respect to  $J$ . The *Jaccard index at instance level* ( $J_i$ ) is then defined as the mean best overlap for all the ground-truth instances in the database, also known as Best Spatial Support score (BSS) [4] or Average Best Overlap (ABO) [12].

Computing the mean of the best overlap on all objects, as done by  $J_i$ , hides the distribution of quality among different objects. As an example,  $J_i = 0.5$  can mean that the algorithm covers half the objects perfectly and completely misses the other half, or can also mean that all the objects are covered exactly at  $J = 0.5$ . This information might be useful to decide which algorithm to use. Computing a histogram of the best overlap would provide very rich information, but then the resulting plot would be 3D (number of proposals, bins, and bin counts). Alternatively, we propose to plot different percentiles of the histogram.

Interestingly, a certain percentile of the histogram of best overlaps consists in computing the number of objects whose best overlap is above a certain Jaccard threshold, which can be interpreted as the best achievable recall of the technique over a certain threshold. We compute the recall at three different thresholds:  $J = 0.5$ ,  $J = 0.7$ , and  $J = 0.85$ .

**Learning Strategy Evaluation:** We first estimate the loss in performance due to not sweeping all the possible values of  $\{N_1, \dots, N_R\}$  in the combination of proposal lists via the proposed greedy strategy. To do so, we will compare this strategy with the full combination on a reduced problem to make the latter feasible. Specifically, we combine the 4 ranked lists coming from the singletons at all scales, instead of the full 16 lists coming from singletons, pairs, triplets, and 4-tuples. We also limit the search to 20000 proposals, further speeding the process up.

In this situation, the mean loss in achievable quality along the full curve of parameterization is  $J_i = 0.0002$ , with a maximum loss of  $J_i = 0.004$  (0.74%). In exchange, our proposed learning strategy on the full 16 ranked lists takes about 20 seconds to compute on the training set of SegVOC12, while the singleton-limited full combination takes 4 days (the full combination would take months).

**Combinatorial Grouping:** We now evaluate the Pareto front optimization strategy in the training set of SegVOC12. As before, we extract the lists of proposals from the three scales and the multiscale hierarchy, for singletons, pairs, triplets, and 4-tuples of regions, leading to 16 lists, ranked by the minimum UCM strength of the regions forming each proposal.

Figure 8 shows the Pareto front evolution of  $J_i$  with respect to the number of proposals for up to 1, 2, 3, and 4 regions per proposal (4, 8, 12, and 16 lists, respectively) at training and testing time on SegVOC12. As baselines, we plot the raw singletons from MCG-UCM-Our, gPb-UCM, and Quadtree; as well as the uniform combination of scales.

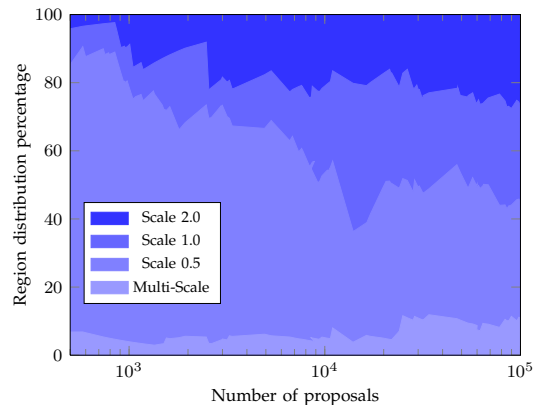


Fig. 10. Region distribution learnt by the Pareto front optimization on SegVOC12.

The improvement of considering the combination of all 1-region proposals (—) from the 3 scales and the MCG-UCM-Our with respect to the raw MCG-UCM-Our (---) is significant, which corroborates the gain in diversity obtained from hierarchies at different scales. In turn, the addition of 2- and 3-region proposals (— and —) noticeably improves the achievable quality. This shows that hierarchies do not get full objects in single regions, which makes sense given that they are built using low-level features only. The improvement when adding 4-tuples (—) is marginal at the number of proposals we are considering. When analyzing the equal distribution of proposals from the four scales (---), we see that the less proposals we consider, the more relevant the Pareto optimization becomes. At the selected working point, the gain of the Pareto optimization is 2 points.

Figure 10 shows the distribution of proposals from each of the scales combined in the Pareto front. We see that the coarse scale (0.5) is the most *picked* at low number of proposals, and the rest come into play when increasing their number, since one can afford more detailed proposals. The multi-scale hierarchy is the one with less weight, since it is created from the other three.

**Pareto selection and ranking:** Back to Figure 8, the red asterisk (\*) marks the selected configuration  $\{\bar{N}_1, \dots, \bar{N}_R\}$  in the Pareto front (black triangle in Figure 7), which is selected at a practical level of proposals. The red plus sign (+) represents the set of proposals after removing those duplicate proposals whose overlap leads to a Jaccard higher than 0.95. The proposals at this point are the ones that are ranked by the learnt regressor (—).

At test time (right-hand plot), we directly combine the learnt  $\{\bar{N}_1, \dots, \bar{N}_R\}$  proposals from each ranked list. Note that the Pareto optimization does not overfit, given the similar result in the training and validation datasets. We then remove duplicates and rank the results. In this case, note the difference between the regressed result in the training and validation sets, which reflects overfitting, but despite this we found it beneficial with respect to the non-regressed result.

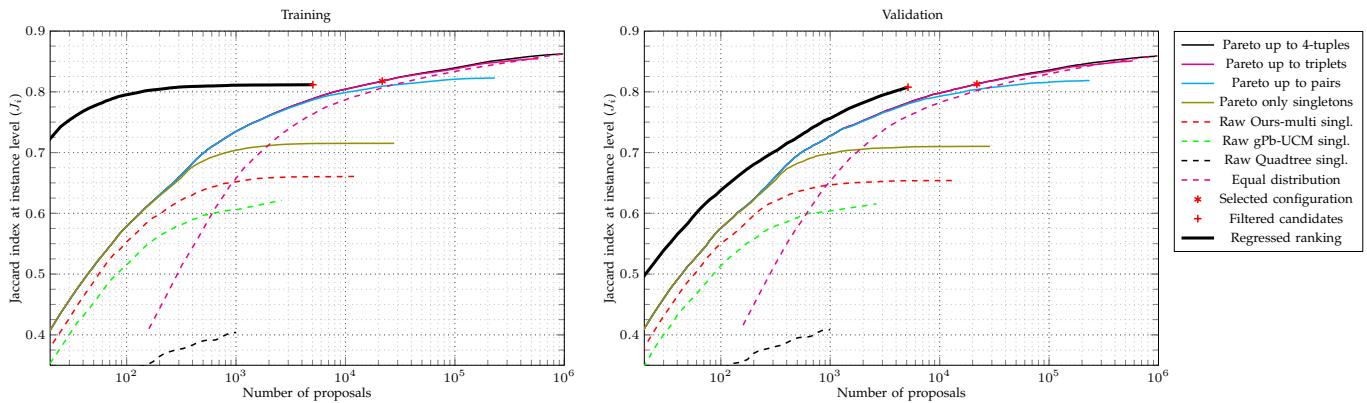


Fig. 8. **Pareto front evaluation.** Achievable quality of our proposals for singletons, pairs, triplets, and 4-tuples; and the raw proposals from the hierarchies on PASCAL SegVOC12 training (left) and validation (right) sets.

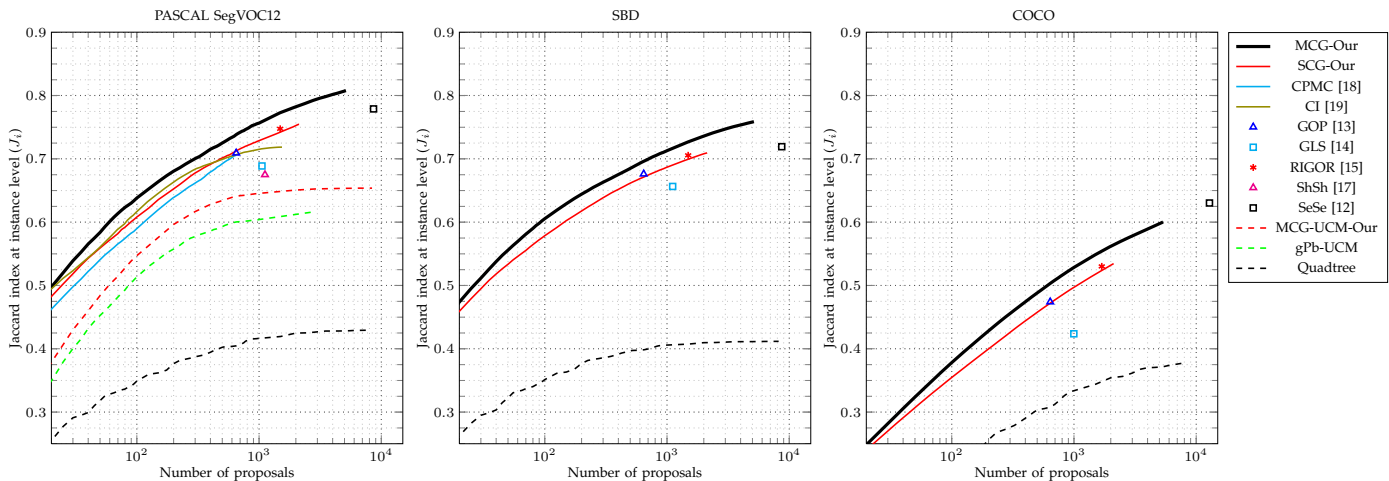


Fig. 9. **Object Proposals: Jaccard index at instance level.** Results on SegVOC12, SBD, and COCO.

In the validation set of SegVOC12, the full set of proposals (i.e., combining the full 16 lists) would contain millions of proposals per image. The multiscale combinatorial grouping allows us to reduce the number of proposals to 5 086 with a very high achievable  $J_i$  of 0.81 (+). The regressed ranking (—) allows us to further reduce the number of proposals below this point.

**Segmented Proposals: Comparison with State of the Art:** We first compare our results against those methods that produce segmented object proposals [13], [14], [15], [12], [16], [17], [18], [19], using the implementations from the respective authors. We train MCG on the training set of SegVOC12, and we use the learnt parameters on the validation sets of SegVOC12, SBD, and COCO.

Figure 9 shows the achievable quality at instance level ( $J_i$ ) of all methods on the validation set of SegVOC12, SBD, and COCO. We plot the raw regions of MCG-UCM-Ours, gPb-UCM, and QuadTree as baselines where available. We also evaluate a faster single-scale version of MCG (*Single-scale Combinatorial Grouping* - SCG), which takes the hierarchy at the native scale only and combines up to 4 regions per proposal. This approach decreases the computational load one order of magnitude while

keeping competitive results.

MCG proposals (—) significantly outperform the state-of-the-art at all regimes. The bigger the database is, the better MCG results are with respect to the rest, which shows that our techniques better generalize to unseen images (recall that MCG is trained only in SegVOC12).

As commented on the measures description,  $J_i$  shows mean aggregate results, so they can mask the distribution of quality among objects in the database. Figure 11 shows the recall at three different Jaccard levels. First, these plots further highlight how challenging COCO is, since we observe a significant drop in performance, more pronounced than when measured by  $J_i$  and  $J_c$ . Another interesting result comes from observing the evolution of the plots for the three different Jaccard values. Take for instance the performance of GOP (▲) against MCG-Ours (—) in SBD. While for  $J=0.5$  GOP slightly outperforms MCG, the higher the threshold, the better MCG. Overall, MCG has specially good results at higher  $J$  values. In other words, if one looks for proposals of very high accuracy, MCG is the method with highest recall, at all regimes and in all databases. In all measures and databases, SCG (—) obtains very competitive results, especially if we take into account

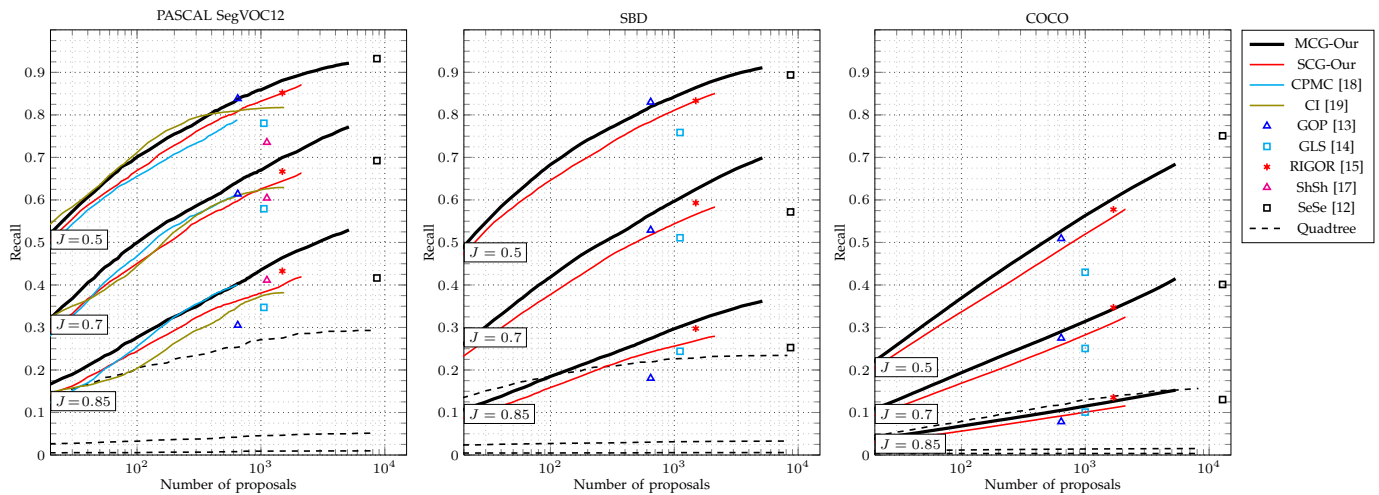


Fig. 11. **Segmented Object Proposals: Recall at different Jaccard levels.** Percentage of annotated objects for which there is a proposal whose overlap with the segmented ground-truth shapes (not boxes) is above  $J = 0.5$ ,  $J = 0.7$ , and  $J = 0.85$ , for different number of proposals per image. Results on SegVOC12, SBD, and COCO.

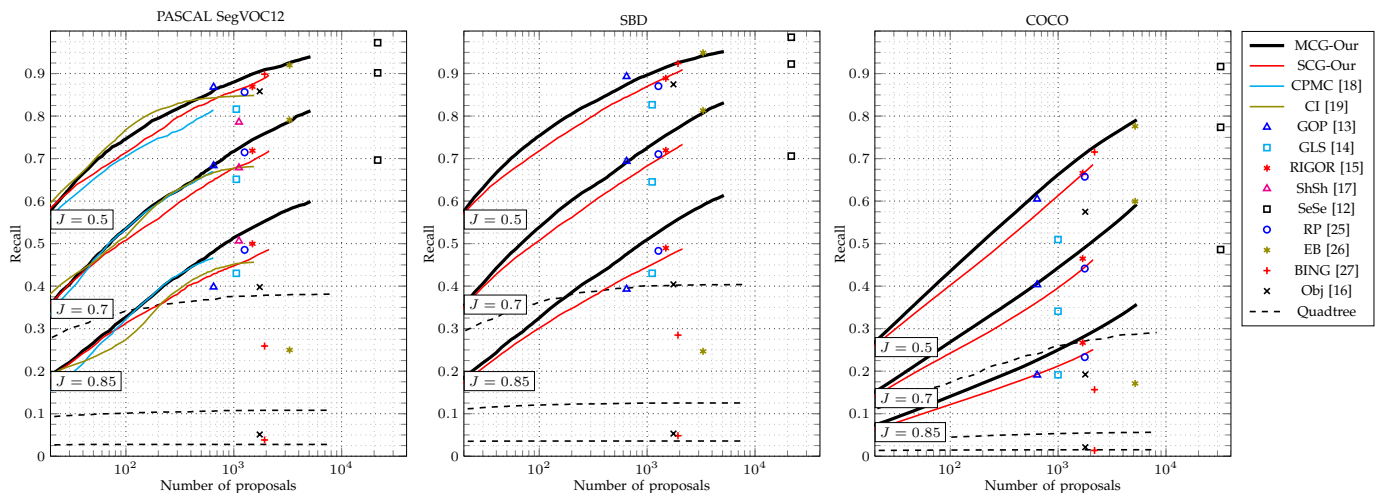


Fig. 12. **Bounding-Box Proposals: Recall at different Jaccard levels.** Percentage of annotated objects for which there is a bounding box proposal whose overlap with the ground-truth boxes is above  $J = 0.5$ ,  $J = 0.7$ , and  $J = 0.85$ , for different number of proposals per image. Results on SegVOC12, SBD, and COCO.

that it is  $7\times$  faster than MCG, as we will see in next sections.

The complementarity of MCG with respect to other proposal techniques, and their combination using the Pareto front optimization is studied in [44].

**Boxes Proposals: Comparison with State of the Art:** Although focused in providing segmented object proposals, MCG may also be used to provide boxes proposals, by taking the bounding box of the segmented proposals and deduplicating them. We add the state of the art in boxes proposals [26], [27], [25], and [16] to the comparison. Figure 12 shows the recall values of the boxes results when compared to the bounding boxes of the annotated objects, for three different Jaccard thresholds.

While many of the techniques specifically tailored to boxes proposals are competitive at  $J = 0.5$ , their perfor-

mance drops significantly at higher Jaccard thresholds. Despite being tailored to segmented proposals, MCG clearly outperforms the state of the art if you look for precise localization of the bounding boxes. Again, SCG is very competitive, especially taking its speed into account.

**MCG and SCG Time per Image:** Table 3 shows the time per image of the main steps of our approach, from the raw image to the contours, the hierarchical segmentation, and the proposal generation. All times are computed using a single core on a Linux machine. Our full MCG takes about 25 s. per image to compute the multiscale hierarchies, and 17 s. to generate and rank the 5038 proposals on the validation set of SegVOC12. Our single-scale SCG produces a segmentation hierarchy of better quality than gPb-ucm [20] in less than 3 seconds and with significant less memory footprint.

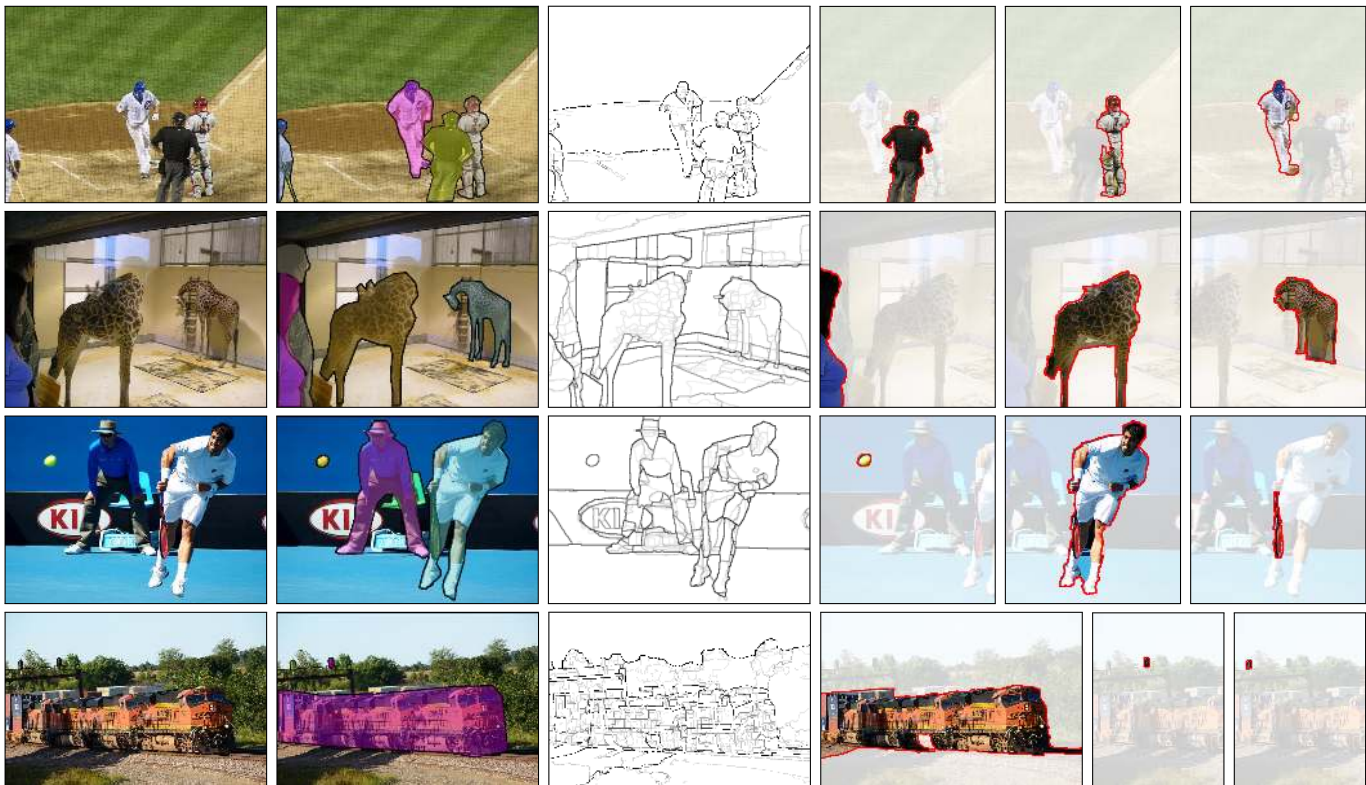


Fig. 13. **COCO Qualitative Results:** Image, ground truth, multi-scale UCM and best MCG proposals among the 500 best ranked. (More qualitative examples in the supplemental material.)

	Contour Detection	Hierarchical Segmentation	Candidate Generation	Total
MCG	$4.6 \pm 1.3$	$20.5 \pm 5.6$	$17.0 \pm 9.8$	$42.2 \pm 14.8$
SCG	$1.0 \pm 0.3$	$2.7 \pm 0.7$	$2.6 \pm 0.4$	$6.2 \pm 1.1$

TABLE 3  
Time in seconds per image of MCG and SCG

Table 4 shows the time-per-image results compared to the rest of state of the art in segmented proposals generation, all run on the same single-core Linux machine.

	Proposal Generation
MCG (Our)	$42.2 \pm 14.8$
SCG (Our)	$6.2 \pm 1.1$
GOP [13]	$1.0 \pm 0.3$
GLS [14]	$7.9 \pm 1.7$
SeSe [12]	$15.9 \pm 5.2$
RIGOR [15]	$31.6 \pm 16.0$
CPMC [18]	$\geq 120$
CI [19]	$\geq 120$
ShSh [17]	$\geq 120$

TABLE 4  
Time comparison for all considered state-of-the-art techniques that produce segmented proposals. All run on the same single-core Linux machine.

**Practical considerations:** One of the key aspects of object proposals is the size of the pool they generate. Depending on the application, one may need more precision and thus a bigger pool, or one might need speed and thus a small pool in exchange for some loss of quality. MCG and SCG provide a *ranked* set of around 5000 and 2000 proposals, respectively, and one can take the  $N$  first in case the specific application needs a smaller pool. From a practical point of view, this means that one does not need to re-parameterize them for the specific needs of a certain application.

In contrast, the techniques that do not provide a ranking of the proposals, need to be re-parameterized to adapt them to a different number of proposals, which is not desirable in practice.

On top of that, the results show that MCG and SCG have outstanding generalization power to unseen images (recall that the results for SBD and COCO have been obtained using the learnt parameters on SegVOC12), meaning that MCG and SCG offer the best chance to obtain competitive results in an unseen database without need to re-train.

Figure 13 shows some qualitative results on COCO.

## 7 CONCLUSIONS

We proposed Multiscale Combinatorial Grouping (MCG), a unified approach for bottom-up segmentation and object proposal generation. Our approach produces

state-of-the-art contours, hierarchical regions, and object proposals. At its core are a fast eigenvector computation for normalized-cut segmentation and an efficient algorithm for combinatorial merging of hierarchical regions. We also present Single-scale Combinatorial Grouping (SCG), a speeded up version of our technique that produces competitive results in under five seconds per image.

We perform an extensive validation in BSDS500, SegVOC12, SBD, and COCO, showing the quality, robustness and scalability of MCG. Recently, an independent study [45], [46] provided further evidence to the interest of MCG among the current state-of-the-art in object proposal generation. Moreover, our object candidates have already been employed as integral part of high performing recognition systems [47].

In order to promote reproducible research on perceptual grouping, all the resources of this project – code, pre-computed results, and evaluation protocols – are publicly available<sup>2</sup>.

**Acknowledgements:** The last iterations of this work have been done while Jordi Pont-Tuset has been at Prof. Luc Van Gool’s Computer Vision Lab (CVL) of ETHZ, Switzerland. This work has been partly developed in the framework of the project BIGGRAPH-TEC2013-43935-R and the FPU grant AP2008-01164; financed by the *Spanish Ministerio de Economía y Competitividad*, and the European Regional Development Fund (ERDF). This work was partially supported by ONR MURI N000141010933.

## REFERENCES

- [1] P. Viola and M. Jones, “Robust real-time face detection,” *IJCV*, vol. 57, no. 2, 2004.
- [2] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection,” in *CVPR*, 2005.
- [3] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part based models,” *TPAMI*, vol. 32, no. 9, 2010.
- [4] T. Malisiewicz and A. A. Efros, “Improving spatial support for objects via multiple segmentations,” in *BMVC*, 2007.
- [5] C. Gu, J. Lim, P. Arbeláez, and J. Malik, “Recognition using regions,” in *CVPR*, 2009.
- [6] J. Carreira, F. Li, and C. Sminchisescu, “Object recognition by sequential figure-ground ranking,” *IJCV*, vol. 98, no. 3, pp. 243–262, 2012.
- [7] A. Ion, J. Carreira, and C. Sminchisescu, “Probabilistic joint image segmentation and labeling by figure-ground composition,” *IJCV*, vol. 107, no. 1, pp. 40–57, 2014.
- [8] P. Arbeláez, B. Hariharan, C. Gu, S. Gupta, L. Bourdev, and J. Malik, “Semantic segmentation using regions and parts,” in *CVPR*, 2012.
- [9] J. Carreira, R. Caseiro, J. Batista, and C. Sminchisescu, “Semantic segmentation with second-order pooling,” in *ECCV*, 2012.
- [10] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *CVPR*, 2014.
- [11] N. Zhang, J. Donahue, R. Girshick, and T. Darrell, “Part-based r-cnns for fine-grained category detection,” in *ECCV*, 2014.
- [12] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders, “Selective search for object recognition,” *IJCV*, vol. 104, no. 2, pp. 154–171, 2013.
- [13] P. Krähenbühl and V. Koltun, “Geodesic object proposals,” in *ECCV*, 2014.
- [14] P. Rantalankila, J. Kannala, and E. Rahtu, “Generating object segmentation proposals using global and local search,” in *CVPR*, 2014.
- [15] A. Humayun, F. Li, and J. M. Rehg, “RIGOR: Recycling Inference in Graph Cuts for generating Object Regions,” in *CVPR*, 2014.
- [16] B. Alexe, T. Deselaers, and V. Ferrari, “Measuring the objectness of image windows,” *TPAMI*, vol. 34, pp. 2189–2202, 2012.
- [17] J. Kim and K. Grauman, “Shape sharing for object segmentation,” in *ECCV*, 2012.
- [18] J. Carreira and C. Sminchisescu, “CPMC: Automatic object segmentation using constrained parametric min-cuts,” *TPAMI*, vol. 34, no. 7, pp. 1312–1328, 2012.
- [19] I. Endres and D. Hoiem, “Category-independent object proposals with diverse ranking,” *TPAMI*, vol. 36, no. 2, pp. 222–234, 2014.
- [20] P. Arbeláez, M. Maire, C. C. Fowlkes, and J. Malik, “Contour detection and hierarchical image segmentation,” *TPAMI*, vol. 33, no. 5, pp. 898–916, 2011.
- [21] P. F. Felzenszwalb and D. P. Huttenlocher, “Efficient graph-based image segmentation,” *IJCV*, vol. 59, p. 2004, 2004.
- [22] X. Ren and L. Bo, “Discriminatively trained sparse code gradients for contour detection,” in *NIPS*, 2012.
- [23] C. J. Taylor, “Towards fast and accurate segmentation,” *CVPR*, 2013.
- [24] M. Maire and S. X. Yu, “Progressive multigrid eigensolvers for multiscale spectral segmentation,” *ICCV*, 2013.
- [25] S. Manén, M. Guillaumin, and L. Van Gool, “Prime Object Proposals with Randomized Prim’s Algorithm,” in *ICCV*, 2013.
- [26] C. L. Zitnick and P. Dollár, “Edge boxes: Locating object proposals from edges,” in *ECCV*, 2014.
- [27] M.-M. Cheng, Z. Zhang, W.-Y. Lin, and P. H. S. Torr, “BING: Binarized normed gradients for objectness estimation at 300fps,” in *CVPR*, 2014.
- [28] D. Hoiem, A. Efros, and M. Hebert, “Recovering occlusion boundaries from an image,” *IJCV*, vol. 91, no. 3, pp. 328–346, 2011.
- [29] Z. Ren and G. Shakhnarovich, “Image segmentation by cascaded region agglomeration,” in *CVPR*, 2013.
- [30] D. Weiss and B. Taskar, “Scalpel: Segmentation cascades with localized priors and efficient learning,” in *CVPR*, 2013.
- [31] L. Najman and M. Schmitt, “Geodesic saliency of watershed contours and hierarchical segmentation,” *TPAMI*, vol. 18, no. 12, pp. 1163–1173, 1996.
- [32] P. Arbeláez, “Boundary extraction in natural images using ultrametric contour maps,” in *POCV*, June 2006.
- [33] D. Martin, C. Fowlkes, and J. Malik, “Learning to detect natural image boundaries using local brightness, color and texture cues,” *TPAMI*, vol. 26, no. 5, pp. 530–549, 2004.
- [34] P. Dollár and C. Zitnick, “Structured forests for fast edge detection,” *ICCV*, 2013.
- [35] <http://www.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/resources.html>.
- [36] J. Pont-Tuset and F. Marques, “Supervised evaluation of image segmentation and object proposal techniques,” *PAMI*, 2015.
- [37] —, “Measures and meta-measures for the supervised evaluation of image segmentation,” in *CVPR*, 2013.
- [38] M. Maire, S. X. Yu, and P. Perona, “Hierarchical scene annotation,” in *BMVC*, 2013.
- [39] M. Everingham, H. Muller, and B. Thomas, “Evaluating image segmentation algorithms using the pareto front,” in *ECCV*, 2006.
- [40] M. Ehrgott, *Multicriteria optimization*. Springer, 2005.
- [41] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results,” <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>.
- [42] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik, “Semantic contours from inverse detectors,” in *ICCV*, 2011.
- [43] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick, “Microsoft COCO: Common Objects in Context,” in *ECCV*, 2014.
- [44] J. Pont-Tuset and L. V. Gool, “Boosting object proposals: From pascal to COCO,” in *ICCV*, 2015.
- [45] J. Hosang, R. Benenson, and B. Schiele, “How good are detection proposals, really?” in *BMVC*, 2014.
- [46] J. Hosang, R. Benenson, P. Dollár, and B. Schiele, “What makes for effective detection proposals?” *PAMI*, 2015.

2. [www.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/mcg/](http://www.eecs.berkeley.edu/Research/Projects/CS/vision/grouping/mcg/)

[47] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik, "Simultaneous detection and segmentation," in *ECCV*, 2014.



**Jordi Pont-Tuset** is a post-doctoral researcher at ETHZ, Switzerland, in Prof. Luc Van Gool's computer vision group (2015). He received the degree in Mathematics in 2008, the degree in Electrical Engineering in 2008, the M.Sc. in Research on Information and Communication Technologies in 2010, and the Ph.D with honors in 2014; all from the Universitat Politècnica de Catalunya, BarcelonaTech (UPC). He worked at Disney Research, Zürich (2014).



**Pablo Arbeláez** received a PhD with honors in Applied Mathematics from the Université Paris-Dauphine in 2005. He was a Research Scientist with the Computer Vision Group at UC Berkeley from 2007 to 2014. He currently holds a faculty position at Universidad de los Andes in Colombia. His research interests are in computer vision, where he has worked on a number of problems, including perceptual grouping, object recognition and the analysis of biomedical images.



**Jonathan T. Barron** is a senior research scientist at Google, working on computer vision and computational photography. He received a PhD in Computer Science from the University of California, Berkeley in 2013, where he was advised by Jitendra Malik, and he received a Honours BSc in Computer Science from the University of Toronto in 2007. His research interests include computer vision, machine learning, computational photography, shape reconstruction, and biological image analysis. He received a National Science Foundation Graduate Research Fellowship in 2009, and the C.V. Ramamoorthy Distinguished Research Award in 2013.



**Ferran Marques** received the degree in Electrical Engineering and the Ph.D. from the Universitat Politècnica de Catalunya, BarcelonaTech (UPC), where he is currently Professor at the department of Signal Theory and Communications. In the term 2002-2004, he served as President of the European Association for Signal Processing (EURASIP). He has served as Associate Editor of the IEEE Transactions on Image Processing (2009-2012) and as Area Editor for Signal Processing: Image Communication, Elsevier (2010-2014). In 2011, he received the Jaume Vives distinction for University Teaching Quality. Currently, he serves as Dean of the Electrical Engineering School (ETSETB-TelecomBCN) at UPC. He has published over 150 conference and journal papers, 2 books, and holds 4 international patents.



**Jitendra Malik** is Arthur J. Chick Professor in the Department of Electrical Engineering and Computer Science at the University of California at Berkeley, where he also holds appointments in vision science and cognitive science. He received the PhD degree in Computer Science from Stanford University in 1985. In January 1986, he joined UC Berkeley as a faculty member in the EECS department where he served as Chair of the Computer Science Division during 2002-2006, and of the Department of EECS during 2004-2006. Jitendra Malik's group has worked on computer vision, computational modeling of biological vision, computer graphics and machine learning. Several well-known concepts and algorithms arose in this work, such as anisotropic diffusion, normalized cuts, high dynamic range imaging, shape contexts and poselets. According to Google Scholar, ten of his papers have received more than a thousand citations each. He has graduated 33 PhD students. Jitendra was awarded the Longuet-Higgins Award for "A Contribution that has Stood the Test of Time" twice, in 2007 and 2008. He is a Fellow of the IEEE and the ACM, a member of the National Academy of Engineering, and a fellow of the American Academy of Arts and Sciences. He received the PAMI Distinguished Researcher Award in computer vision in 2013 and the K.S. Fu prize in 2014.