

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Scale-Aware Alignment of Hierarchical Image Segmentation

Anonymous CVPR submission

Paper ID 752

Abstract

Image segmentation is a key component in many computer vision systems, and it is recovering a prominent spot in the literature as methods improve and overcome their limitations. The outputs of most recent algorithms are in the form of a hierarchical segmentation, which provides segmentation at different scales in a single tree-like structure. Commonly, these hierarchical methods start from some low-level features, and are not aware of the scale information of the different regions in them. As such, one might need to work on many different levels of the hierarchy to find the objects in the scene. This work tries to modify the existing hierarchical algorithm by improving their alignment, that is, by trying to modify the depth of the regions in the tree to better couple depth and scale. To do so, we first train a regressor to predict the scale of regions using mid-level features. We then define the anchor slice as the set of regions that better balance between over-segmentation and under-segmentation. The output of our method is an improved hierarchy, re-aligned by the anchor slice. To demonstrate the power of our method, we perform comprehensive experiments on the BSDS500 dataset, which shows that our method, as a post-processing step, can significantly improve the quality of the hierarchical segmentation representations. We also prove that the improvement generalizes well across different algorithms, with a low computational cost.

1. Introduction

Generic image segmentation has been part of computer vision and image processing communities since the advent of these fields many decades ago. The definition of the problem, although vague, is easy to give and understand: “to divide the pixels of an image into different pieces, where each piece represents a distinguished *thing* in the image.” Martin *et al.* [16] provided these instructions to annotators to create the Berkeley Segmentation Database (BSDS), which proved that the problem of image segmentation was, indeed, well defined, as humans provided consistent partitions of the images *up to refinement*. In other words, image

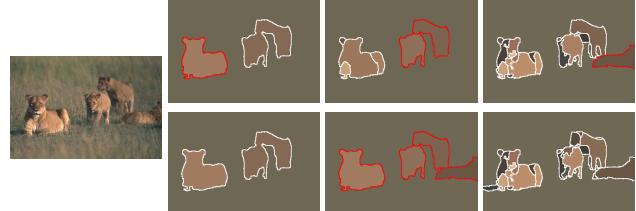


Figure 1. Example of improved hierarchy alignment: The original hierarchy (top row) needs three different flat partitions to represent the four objects (highlighted in red). Our aligned hierarchy (bottom row) correctly puts all objects in the same level.

segmentation is inherently a multi-scale problem.

We refer to *flat* image segmentation techniques as those whose output is a single partition of the image pixels into sets [25, 4, 9]. In these cases, in order to capture the aforementioned multi-scale nature of objects, one needs to sweep different parameterizations to obtain multiple partitions that contain the different scales when working with flat segmentation techniques.

On the other hand, hierarchical segmentation produces a single multi-scale structure that aims at capturing the objects at all scales [1, 13, 24, 21, 2]. These types of structures have been successfully used in image filtering [24], semantic segmentation [14], object proposals generation [2], or video segmentation [27, 26].

The representation power of these hierarchies comes at a cost, however, which is the difficulty to handle them from a practical (coding) point of view. While a flat partition can be represented by a matrix of labels of each pixel, hierarchical structures need a much more complex representation. In this context, the Ultrametric Contour Map (UCM) [1] representation is the one that gained more traction and it is widely used in the literature. In it, *flattening* the hierarchy can be achieved by simply *thresholding* the UCM.

The process of *flattening* or *pruning* a hierarchy is therefore of paramount importance for segmentation, because it is the main proxy used towards the final application. This work presents a novel technique to improve the flattening of any given hierarchy, that is, to get better flat partitions from the same hierarchical segmentation.

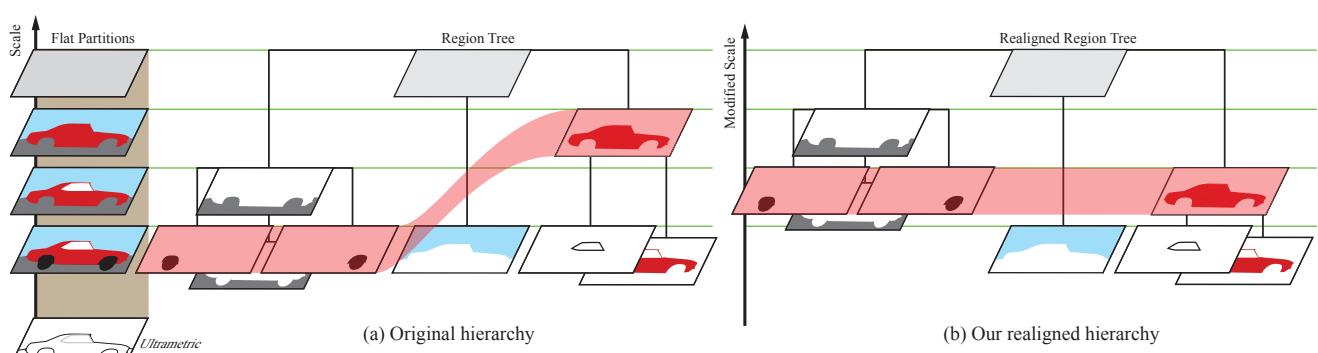


Figure 2. Our proposed hierarchy realignment: Given a hierarchy (a) in which the objects at the same scale are not well aligned (represented in the same scale level), we produce a realigned hierarchy (b) that has the similar-scale regions in the same level.

Figure 1 motivates this work. In the first row we can see different flat partitions extracted from the same hierarchy. To get the regions representing the four lions we need to search in three different flat partitions, extracted at three different levels of the hierarchy. The second row shows our results, where the same hierarchy is *aligned* to have all objects represented in the same flat partition.

In other words, the threshold level of the hierarchy better relates with the scale of the objects, not only in the same image, but also across images. To further grasp the intuition of our work, Figure 2 shows a UCM and its interpretation as a region tree (a). In it, the needed regions to form the car are spread into different scale levels (thresholds of the UCM), as marked by the red band. Our proposed realigned hierarchy (b) aims at containing them all in the same scale.

Since the hierarchies are constructed based on low-level features (edges or color), the scale of the objects is not imposed to be coherent. We propose to learn the concept of object scale from mid-level features within the hierarchy. Our objective is to take advantage of these mid-level features as much as possible without getting to high-level features that would allow us to go beyond scale. This way, the global approach would be to construct the hierarchies using low-level features, and then exploit mid-level features to realign them, thus taking the maximum advantage of the most simple features possible.

Our alignment also aims at providing a global alignment among different images, that is, providing levels of scale that keep meaning even when changing images, allowing higher-level methods to generalize in a more straightforward manner. Specifically, we train a regressor to predict whether each region of the hierarchy is oversegmented, undersegmented, or correctly segmented; and we rescale the hierarchy according to the prediction of this classifier. Back to the example in Figure 1, the majority of regions in the first column (bottom) are undersegmented, in the middle column they are correctly segmented, and oversegmented in the last column.

We perform comprehensive experiments on BSDS500 using four different hierarchical segmenters. We obtain consistent improvements on all hierarchies which proves the usefulness of our approach and its generalization power. The remainder of the paper is organized as follows. First, Section 2 gives a brief overview of the related work. Then Section 3 presents our algorithm for re-scaling and aligning hierarchies. We demonstrate the effectiveness of our method in the experiments in Section 4 and draw the conclusions in Section 5.

2. Related Work

Hierarchical Segmentation There is a rich literature about hierarchical segmentation. As stated in the introduction, our main focus in this paper is not to develop a better hierarchical segmentation algorithm, but to provide a better alignment of a given hierarchy. Hierarchical segmentation typically starts from various local information embedded in an affinity matrix, such as Pointwise Mutual Information [12], or multiscale local brightness, color, and texture cues [1]. It then greedily constructs a hierarchy of regions by iteratively merging the most similar sets of regions according to a certain metric. The result of hierarchical segmentation is commonly represented as an Ultrametric Contour Map (UCM), where different levels of segmentation can be produced by applying different thresholds to UCM. This work proposes to realign the hierarchies in order to make the thresholds of the UCM more closely related to the scale of objects. Hierarchical segmentation has become the major trend in image segmentation and most of top-performance segmenters [1, 2, 21, 13] fall into this category.

Multiple Segmentations Working with multiple segmentations at the same time has been used in the computer vision community for a long time, with the idea that, while none of the segmentations is likely to partition the image

perfectly, some parts in some segmentations might be useful. Hoiem *et al.* [10] use this idea to estimate the scene structure. A similar idea was exploited by Russell *et al.* [23] to discover objects, and by Malisiewicz *et al.* [15] to improve the spatial support of regions for recognition. By re-aligning the hierarchies we aim to minimize the number of partitions from a hierarchy needed to obtain reasonable results, since we concentrate same-scale regions in the same partition.

Our work also shares some similarities with [27], where they flatten supervoxel hierarchies by finding a slice with uniform entropy. However, they formulate the problem as a quadratic integral programming, which makes the optimization NP-hard. On the contrary, we use a simple yet effective dynamic programming algorithm to find a globally-optimal slice.

Predicting Segmentation Quality by Classification

Classification has been exploited to predict segmentation quality in many works. Ren *et al.* [20] use a linear classifier base on Gestalt features [17] to distinguish good and bad segmentations. Their negative training data are generated by randomly placing a ground-truth mask over an image. A similar idea is used to select parameters by Peng *et al.* [18] to select λ in graph-cut based interactive segmentation. They compute the segmentation with different λ , then select the one with highest predicted quality. More recently, Carreira *et al.* [3], Arbelaez *et al.* [2], and Endres *et al.* [7] use a regression forest to predict the good overlap between segments (object proposals) and ground truth objects. We use similar features to [3], which are based on graph partition properties, region properties, and Gestalt properties.

3. Flattening and Re-scaling Hierarchies

As discussed in the introduction, while segmentation hierarchies contain a rich multiscale decomposition of the image, it is not trivial to distill such knowledge because the hierarchies generated by current methods are not fully scale-aware. Simply taking a layer yields a segmentation of which some parts are under-segmented while others are over-segmented. In this section, we present our method which aligns the scales of segmentation hierarchies, making image hierarchies easier to use in practice. We start with scale labeling, and then present the alignment strategy.

3.1. Flattening Hierarchies via Scale Labeling

Let's denote by \mathcal{T} the segmentation tree of image I , with node v_i indicating its i th node. The nodes correspond to regions (segments) of I . Given \mathcal{T} , our task is to find a tree slice \mathcal{L} to segment all nodes v_i 's (segments) into three groups: \mathcal{L}^- , \mathcal{L} , and \mathcal{L}^+ indicating under-, properly- and over-segmented, respectively. See Figure 3(b) for an example of nodes in the three groups.

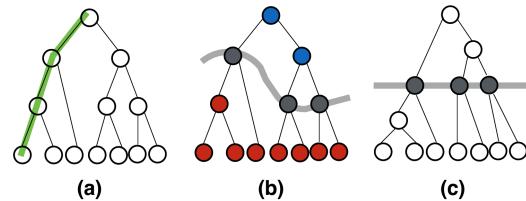


Figure 3. Examples of the slices and paths of the segmentation tree, where one path of the tree is shown in green (a) and one slice is shown in grey (b). In (b), all nodes in blue are in \mathcal{L}^- , and all nodes in red are in \mathcal{L}^+ . Our approach re-aligns the hierarchy using the anchor slice. The aligned tree is shown in (c).

The visual representation of a slice can be seen in Figure 2 as red bands covering different regions. An example of the three types of slices can be found in Figure 1 (bottom row), where the left partition is mainly oversegmented, the middle one correctly segmented, and the right one undersegmented.

We formulate the problem as a three-class labeling problem. For each node v_i , we use $x(v_i) \in \{-1, 0, 1\}$ as its class label, with -1 , 0 , and 1 indicating the membership of v_i to \mathcal{L}^- , \mathcal{L} , and \mathcal{L}^+ respectively. Assume now that a function $f(v_i) : v_i \rightarrow [-1, 1]$ is provided to measure the granularity of image segments, where negative values stand for under-segmented, 0 for properly-segmented, and positive for over-segmented regions. The magnitude of $f(v_i)$ signals the deviation from being properly-segmented. Section 3.1.2 present the proposed learning algorithm for $f(v_i)$.

The labeling of all v_i 's could be done by greedily taking the best-scoring class for each node. However, not any labeling represents a valid slice of the tree. Following the definition in [19, 27], a tree slice is a set of nodes such that every path $\mathcal{P}_n, n \in \{1, 2, \dots, N\}$ from the leaf node \bar{v}_n to the root node v_0 contains one and only one node v in the slice. See Figure 3 for the examples of the slices and paths.

From the nature of segmentation hierarchies, the labels of parent nodes v_i^p should be equal or smaller than their child nodes v_i . Intuitively, if a region is correctly segmented, the parent cannot be oversegmented. The same way, the parent of an undersegmented region will also be undersegmented. Putting the two constraints together, the labeling problem can be formulated as:

$$\begin{aligned} \hat{\mathbf{X}} &= \arg \min_{\mathbf{X}} E(\mathbf{X}) \\ E(\mathbf{X}) &= \sum_{v_i \in \mathcal{L}} \#(v_i) \cdot \|f(v_i)\|^2 + \lambda \sum_{v_i \notin \mathcal{L}} \#(v_i) \cdot l(v_i) \\ \text{s.t. } \forall n : \sum_{v \in \mathcal{P}_n} \mathbb{1}_{\mathcal{L}}(v) &= 1 \\ \forall v : x(v) &>= x(v^p) \end{aligned} \quad (1)$$

324 **Algorithm 1** Dynamic Programming in a Tree
325

```

326   Input: tree node  $v_i$ 
327   if  $v_i$  is a leaf node then
328      $C_{v_i} \leftarrow \#(v_i) \cdot \max(0, f(v_i))$ 
329      $E_{v_i}^* \leftarrow \#(v_i) \cdot \|f(v_i)\|^2$ 
330   else
331      $C_{v_i} \leftarrow \sum_{v_j \in \{v^c\}} C_{v_j} + \#(v_i) \cdot \max(0, f(v_i))$ 
332      $E_{v_i}^* \leftarrow \min(\sum_{v_j \in \{v^c\}} E_{v_j}^* + \lambda \cdot \#(v_i) \cdot \max(0, -f(v_i)), \#(v_i) \cdot \|f(v_i)\|^2 + \lambda \cdot \sum_{v_j \in \{v^c\}} C_{v_j})$ 
333   end if
334   return  $C_{v_i}, E_{v_i}^*$ 
335

```

337 where $\#(v)$ is the size (number of pixels) of segment (node)
 338 v , λ is a weighting value for the two energy terms, and $l(v_i)$
 339 is the loss function defined for $v_i \in \{\mathcal{L}^-, \mathcal{L}^+\}$:

$$340 \quad l(v_i) = \max(0, f(v_i) \cdot x(v_i)). \quad (2)$$

341 Intuitively, the loss function penalizes two contradictory
 342 cases: (i) segments in the group of under-segmented that
 343 receive positive scores; and (ii) segments in the group of
 344 over-segmented that receive negative scores. The problem
 345 will be solved via dynamic programming, as explained in
 346 the following section.

347 **3.1.1 Inference by Dynamic Programming**

348 The optimization problem in Equation 1 is highly structured
 349 and can be solved recursively by Dynamic Programming.
 350 For the subtree rooted at node v , its optimal slice $\mathcal{L}(v)$ is
 351 either the node v itself or the union of the optimal slices of
 352 all its child nodes v^c 's, depending on whose energy is lower.
 353 Thus, the problem has optimal substructure [6] and so it
 354 naturally fits to the framework of dynamic programming to
 355 find the global optimal solution.

356 The problem proceeds from bottom to the top of the tree.
 357 For each subtree rooted at the current node v , the energy of
 358 $v \in \mathcal{L}(v)$ is computed and the energy of the optimal slices
 359 of all its child nodes is requested for comparison. The algo-
 360 rithm traverses back, and all comparison will be completed
 361 when the algorithm reaches the root node, and the global
 362 optimal of Equation 1 is obtained. The method is highly ef-
 363 ficient with complexity $\mathcal{O}(N)$, where N is the total number
 364 of nodes. The global optimal of the energy can be found by
 365 applying Algorithm 1 to the root node, and the optimal slice
 366 is the set of nodes labeled to 0.

367 **3.1.2 Predicting the Scales of Segments**

368 In order to predict the scales (under-, properly-, or over-
 369 segmented) of the segments, we follow the route of modern
 370 computer vision systems to learn a predictor from human-
 371 annotated training data. To this end, we define a mea-

372 sure to compare the scale of an image segment r to that
 373 of the corresponding human-annotated segment g . The cor-
 374 respondence is built up by computing the overlap between
 375 computer-generated segments and human-annotated ones
 376 – the most-overlapping human-annotated segment is taken
 377 as the ground-truth for the computer-generated ones. The
 378 overlap is computed with the Intersection over Union (IoU).

379 After having the ground-truth segment g , the scale of the
 380 segment r is then defined as:

$$381 \quad S(r) = \frac{\#(g) - \#(r)}{\max(\#(r), \#(g))}. \quad (3)$$

382 The value of $S(r)$ is in $[-1, 1]$, with negative values
 383 for *under-*, 0 for *properly-* and positive values for *over-*
 384 *segmented* regions, which casts to what we expected from
 385 $f(v)$ (c.f. Section 3.1).

386 With Equation 3, the *scales* of the segments by segmen-
 387 tation methods can be computed and used as the training
 388 data to train our scale predictor.

389 As to the learning method, we employ a regression fore-
 390 est [] as the predictor $f(v)$. As to the features, we use a
 391 set of low-, and middle-level features, mainly following the
 392 work done for object proposals [3, 2]. The features are
 393 designed to capture a variety of region properties such as
 394 Gestalt properties and Graph properties. The detailed list of
 395 the features will be provided in the supplementary material.

396 The main difference between our prediction and the pre-
 397 vious work [3, 20, 2] is that they predict the quality of seg-
 398 ments, while we predict the scale of the segments. Although
 399 numerous measures have been proposed, it is still very hard
 400 to quantify the quality of segments. The granularity of seg-
 401 ments, however, is easier to quantify, and it also provides
 402 more specific information such as under-segmented or over-
 403 segmented.

404 **3.2. Hierarchy Re-scaling with Labeled Scales**

405 After setting the optimal slice, we use it as an anchor
 406 slice to stretch the tree correspondingly. In our experiments,
 407 we use the threshold value of each optimal node as a control
 408 point, and linearly interpolate the original hierarchy.

409 Segmentation tree can be represented in the form of
 410 Ultra-Contour-Map(UCM) [1]. UCM is a matrix with the
 411 size $(2h - 1) * (2w - 1)$, where the h is the height of the
 412 original image, and w is the width. Therefore for each pair
 413 of neighboring pixels in the image, the value in the UCM
 414 matrix represents their boundary strength. And usually it is
 415 a value between 0 and 1. Segmentation of a certain scale can
 416 be extracted by thresholding the UCM with corresponding
 417 value. Without losing any generality, our algorithm directly
 418 manipulate on UCM due to its popularity and simplicity.
 419 The algorithm is summarized in Algorithm 2. We assume
 420 the functions **Boundary** to find the corresponding elements

432 **Algorithm 2** Rescaling Hierarchy

433 **Input:** Optimal Slice \mathcal{S} , UCM map M_{ucm}

434 **for** $r \in \mathcal{S}$ **do**

435 $b \leftarrow \text{Boundary}(r)$

436 $a \leftarrow \text{InnerArea}(r)$

437 $m \leftarrow \min(M_{ucm}(b))$

438 $M_{ucm}(a) \leftarrow M_{ucm}(a) * 0.5/m$

439 **end for**

440 $b_{all} \leftarrow \text{Boundary}(\mathcal{S})$

441 $m_{min} \leftarrow \min(M_{ucm}(b_{all}))$

442 $M_{ucm}(b_{all}) \leftarrow m_{min} + \frac{(M_{ucm}(b_{all}) - m_{min})}{2(1 - m_{min})}$

445
446 of boundary of a region r in the UCM matrix, and **Inner-**
447 **Area** to find its inner area.

448 By the presented algorithm, we align the optimal slice
449 to threshold 0.5, which makes it easier for later use. And
450 no information in the original hierarchy is lost during the
451 rescaling process.

452 **4. Experiments**

453 We evaluate our approach on the segmentation hierar-
454 chies generated by multiple different segmentation meth-
455 ods, and further examine its usefulness on the task of object
456 segmentation. The goal is to demonstrate that the proposed
457 method is able to improve general segmentation hierarchies
458 and the improvement is reflected to high-level vision tasks
459 as well.

460 **4.1. Experiment Settings**

461 **Dataset:** We benchmark the performance of our ap-
462 proach on BSDS500 dataset [1], which includes 500 im-
463 ages, with 200 for training, 100 for validation, and 200 for
464 testing. Each image is annotated by 5 different people on
465 average. As to evaluation, we employ three standard met-
466 rics: Segmentation Covering (SC), Probabilistic Rand In-
467 dex (PRI), and Variation of Information (VI). Readers are
468 referred to [1] for details about the dataset and the evalua-
469 tion metrics.

470 **Candidate methods:** As to the candidate hierarchical
471 segmentation methods, we chose the following four meth-
472 ods due to their popularity and excellent performance:

- 473 • gPb-owt-ucm [1]: a widely-used hierarchical segmen-
474 tation method. Discriminative features are learned for
475 local boundary detection and spectral clustering is ap-
476 plied on top of it for boundary globalization.
- 477 • MCG [2]: a unified framework for segmentation and
478 object proposals. It combines information from multi-
479 ple resolutions of the image and achieves the state-of-
480 the-art results for both image segmentation and object
481 proposals.

- 482 • SCG [2]: the single-resolution, faster version version
483 of MCG. It gets competitive results at a fraction of the
484 cost of MCG.

- 485 • PMI [12]: a recent work for unsupervised boundary
486 detection. It can be applied for image segmentation as
487 well in order to generate a hierarchical segmentation.

488 **Training:** The training set and the validation set of
489 BSDS500 are pooled together as the training set for our re-
490 gression forest. The four segmentation methods are used
491 to generate hierarchies, over which the training samples
492 (segments) are extracted. We trained method-specific and
493 method-agnostic regression forests as the scale predictor.
494 Since a large portion of regions in the hierarchies are very
495 small and features extracted from them are not reliable, we
496 exclude regions smaller than 50 pixels for the training of the
497 predictor.

498 Specifically, for each region r , we find its corresponding
499 ground-truth region g by taking the human-annotated one
500 with the highest covering score. The relative scale of r is
501 then computed with Equation 3 for the regression target of
502 r . As to the features for r , we draw on the success of object
503 proposals [3, 2]. There, a large pool of middle-level features
504 have been defined for segment description. The features
505 used are summarized as follows:

- 506 • Graph partition properties: cut, ratio cut, normalized
507 cut, unbalanced normalized cut.
- 508 • Region properties: area, perimeter, bounding box size,
509 major and minor axis lengths of the equivalent ellipse,
510 eccentricity, orientation, convex area, Euler number.
- 511 • Gestalt properties: inter- and intra-region texton sim-
512 ilarity, inter- and intra-region brightness similarity,
513 inter- and intra-region contour energy, curvilinear con-
514 tinuity, convexity.

515 Readers are referred to [3] for the details of these features.
516 This list is definitely not exhaustive. More high-level fea-
517 tures, e.g. by object detection, could be added.

518 Although these features are simple, extracting them for
519 all layers of the segmentation hierarchies can be costly. Do-
520 ing so is also unnecessary, as most segments at one layer
521 are very similar to those from the parent layer and the child
522 layer. Thus, we only extract the features from a subset of
523 layers sub-sampled from the hierarchies. The layers are uni-
524 formly sampled over the range of UCM values.

525 As to the parameters of our method, we set 100 trees for
526 the random forest. λ in Equation 1 is set to 0.1 to balance
527 information from the three groups, because there are more
528 segments over and under the optimal slice \mathcal{L} .

	Covering (\uparrow)		PRI (\uparrow)		VI (\downarrow)	
	ODS	OIS	ODS	OIS	ODS	OIS
MCG	0.61	0.68	0.83	0.86	1.57	1.39
MCG-aligned	0.63	0.68	0.83	0.86	1.53	1.38
SCG	0.60	0.66	0.83	0.86	1.63	1.42
SCG-aligned	0.61	0.67	0.83	0.86	1.61	1.41
gpb	0.59	0.65	0.83	0.86	1.69	1.48
gpb-aligned	0.60	0.66	0.83	0.86	1.66	1.46
PMI	0.53	0.59	0.76	0.81	2.03	1.80
PMI-aligned	0.54	0.59	0.76	0.81	2.01	1.80

Table 1. The results of our aligned hierarchies with a comparison to the original hierarchies.

4.2. Results

The results of our method evaluated on top of the four candidates segmentation approaches are summarized in Table 1. As shown in the table, the improvements achieved by our alignment are considerable and, more importantly, they are consistent across different methods. The method improves more on ODS than OIS, this is because OIS accesses the ground-truth segmentations to search for the best-performing threshold, which somehow diminish the learned knowledge. We argue that ODS is more practical than OIS in a real vision systems, because for real applications there is no human-annotated segmentations.

In Figure 4 we show segmentation examples of MCG and aligned MCG by our method. As the figure shows, the aligned hierarchies generate characteristics closer to what human expect when flat segmentations are sampled out of the hierarchies. More particularly, after alignment, sampled segmentations of the hierarchies generate consistent responses across all parts of the image: all parts under-segmented, to all parts properly-segmented, and finally to all over-segmented while sampling from the top to the bottom of the hierarchies. This alignment greatly simplifies the use of hierarchical image segmentation for other vision tasks. Figure 6 shows qualitative results with different hierarchies. Our approach shows a consistent improvement over the original results. Since our approach is scale-aware, regions at the same level of the hierarchy are of similar scales across all areas of the images after the alignment. Thus our method demonstrates better ability of preserving regions across images.

We also tested the method in the scenario where the random forests are trained with segments from all of the four methods, and applied to all of them at test time. This gives slightly poorer results but in turn shows that our method can be applied in a method-agnostic approach. The numbers are reported in the supplementary material.

	Covering \uparrow		PRI (\uparrow)		VI (\downarrow)	
	ODS	OIS	ODS	OIS	ODS	OIS
Ncut [25]	0.45	0.53	0.78	0.80	2.23	1.89
Felz-Hutt [9]	0.52	0.57	0.80	0.82	2.21	1.87
Mean Shift [5]	0.54	0.58	0.79	0.81	1.85	1.64
Hoiem [11]	0.56	0.60	0.80	0.77	1.78	1.66
gPb-owt-ucm [1]	0.59	0.65	0.83	0.86	1.69	1.48
ISCRA [22]	0.59	0.66	0.82	0.85	1.60	1.42
PFE+mPb [28]	0.62	0.67	0.84	0.86	1.61	1.43
PFE+MCG [28]	0.62	0.68	0.84	0.87	1.56	1.36
MCG [2]	0.61	0.67	0.83	0.86	1.57	1.39
MCG+Ours	0.63	0.68	0.83	0.86	1.53	1.38

Table 2. Segmentation results on BSDS500 test set, with a comparison to the state-of-the-art competitors.

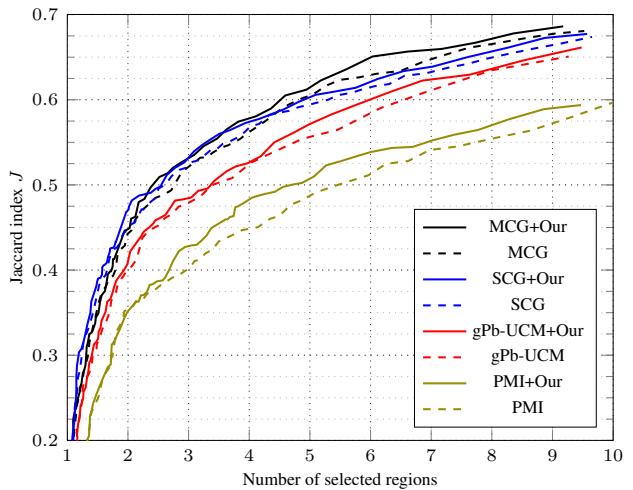


Figure 5. Flattened hierarchies for object detection: Achievable quality by an oracle with respect to the number of regions needed

4.3. Comparison to Other Methods

As the previous section shows, the MCG aligned by our method generally performs the best. Here, we compare MCG-aligned to other competing methods. The results are summarized in Table 2 and demonstrate that segmentation quality can be improved by our alignment. In particular, the aligned MCG achieves the best result in Covering and VI. After alignment, the results are on par with the newest method of PFE+MCG [28]. It is noteworthy that our method and theirs are complementary, and the combination of the two may yield even better results. Their method is to improve feature embedding for a better local distance measure, while we aim to improve the hierarchy of any segmentation method.

4.4. Evaluation towards Object Segmentation

Segmentation *per se* is rarely the final objective of real applications, it is rather a middle tool towards, for instance,



Figure 4. Results of MCG(first row) and MCG results improved by our approach(second row). Original images are shown in the left most. Segmentations of optimal-dataset-sclae(ODS) are given in the middle. And from left to right are different scales, from fine to coarse. Red bounding box indicates the scale with best results achieved by MCG, and blue box for ours. It can be seen that our approach provides better alignment, both across images and within one image.

object segmentation [2] or semantic segmentation [14]. This section is devoted to show that better aligned hierarchies also help in this scenario.

In this case we perform the evaluation using the object

annotations provided on the BSDS300 set by [8] (we retrain on only BSDS300 train instead of BSDS500). The intuitive idea is to measure how well we can segment these objects by *selecting* regions from the different flattened hierarchies.

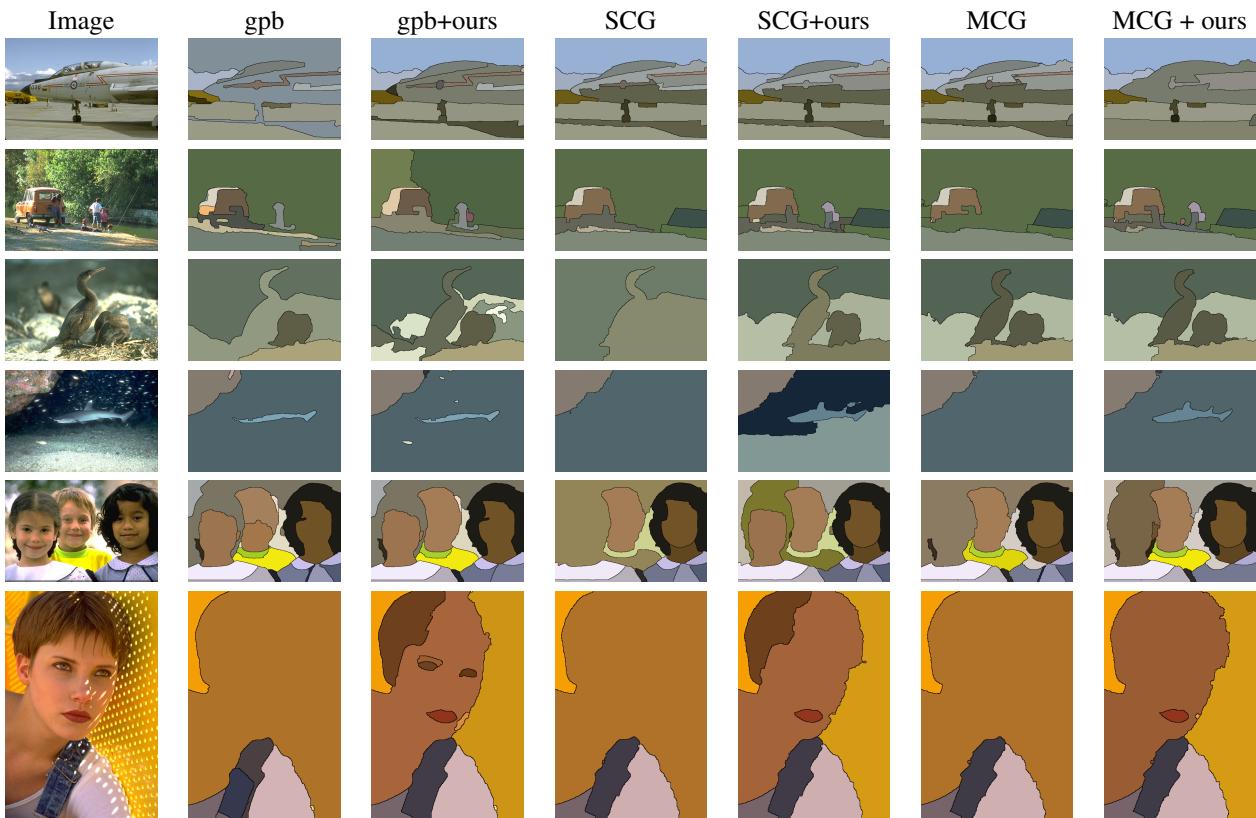


Figure 6. Comparison of segmentation results, hierarchies are flattened by optimal dataset scale(ODS).

Figure 5 shows the achievable quality that an oracle could reach if selecting the regions from the original hierarchies or the ones with our newly-proposed alignment. The X axis corresponds to the number of needed regions, i.e., the lower the better.

We can observe that the aligned hierarchies consistently need less regions to get the same quality in all the tested hierarchies. In PMI, for instance, we need to select 5 regions to achieve the same quality that we can get with 4 on the aligned hierarchy. The combinatorial space of all possible 4-region combinations is significantly smaller and thus the search is more probable to succeed. On the other direction, if we limit the number of regions we get improvements up to 3 points (9%) in the achievable quality.

4.5. Running Time

Our approach takes 3 seconds in total for each image, of which 2.39 seconds are spent on feature extraction from the segments. The prediction of regression forest takes about 0.45 seconds, and the dynamic programming takes 0.05 seconds for the inference. Finally, 0.11 seconds are spent for re-scaling the UCM. All times are measured on a standard desktop machine without parallelization.

5. Conclusion

In this work, we presented a novel technique to align segmentation hierarchies. We proposed a method to learn and predict the scale of segments. We formulated the scale prediction for the segments in a hierarchy as a graph label problem, which is solved by dynamic programming. With the labeled scales as constraints, we then re-align the segmentation hierarchies by stretching the UCM maps. The method is evaluated on four different segmentation hierarchies on BSDS500, and it consistently improves their quality. We also showed that the improvement of segmentation hierarchies by our alignment is reflected well to a higher-level task of getting object segmentations.

References

- [1] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(5):898–916, 2011. [1](#), [2](#), [4](#), [5](#), [6](#)
- [2] P. Arbelaez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 328–335. IEEE, 2014. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#)

- 864 [3] J. Carreira and C. Sminchisescu. Constrained parametric
865 min-cuts for automatic object segmentation. In *Computer
866 Vision and Pattern Recognition (CVPR), 2010 IEEE Conference
867 on*, pages 3241–3248. IEEE, 2010. 3, 4, 5 918
868 [4] D. Comaniciu and P. Meer. Mean shift: a robust approach
869 toward feature space analysis. *IEEE Transactions on Pattern
870 Analysis and Machine Intelligence*, 24(5):603 –619, may
871 2002. 1 919
872 [5] D. Comaniciu and P. Meer. Mean shift: A robust approach
873 toward feature space analysis. *Pattern Analysis and Machine
874 Intelligence, IEEE Transactions on*, 24(5):603–619, 2002. 6 920
875 [6] T. H. Cormen. *Introduction to algorithms*. MIT press, 2009.
876 4 921
877 [7] I. Endres and D. Hoiem. Category-independent object pro-
878 posals with diverse ranking. *Pattern Analysis and Machine
879 Intelligence, IEEE Transactions on*, 36(2):222–234, 2014. 3 922
880 [8] I. Endres and D. Hoiem. Category-independent object pro-
881 posals with diverse ranking. *IEEE TPAMI*, 36(2):222–234,
882 2014. 7 923
883 [9] P. F. Felzenszwalb and D. P. Huttenlocher. Efficient graph-
884 based image segmentation. *International Journal of Com-
885 puter Vision*, 59(2):167–181, 2004. 1, 6 924
886 [10] D. Hoiem, A. A. Efros, and M. Hebert. Geometric context
887 from a single image. In *Computer Vision, 2005. ICCV 2005.
888 Tenth IEEE International Conference on*, volume 1, pages
889 654–661. IEEE, 2005. 3 925
890 [11] D. Hoiem, A. A. Efros, and M. Hebert. Recovering occlusion
891 boundaries from an image. *International Journal of Com-
892 puter Vision*, 91(3):328–346, 2011. 6 926
893 [12] P. Isola, D. Zoran, D. Krishnan, and E. H. Adelson. Crisp
894 boundary detection using pointwise mutual information. In
895 *Computer Vision–ECCV 2014*, pages 799–814. Springer,
896 2014. 2, 5 927
897 [13] T. H. Kim, K. M. Lee, and S. U. Lee. Learning full pairwise
898 affinities for spectral segmentation. *Pattern Analysis and
899 Machine Intelligence, IEEE Transactions on*, 35(7):1690–
900 1703, 2013. 1, 2 928
901 [14] V. Lempitsky, A. Vedaldi, and A. Zisserman. A pylon model
902 for semantic segmentation. In *Advances in Neural Infor-
903 mation Processing Systems (NIPS)*, 2011. 1, 7 929
904 [15] T. Malisiewicz and A. A. Efros. Improving spatial support
905 for objects via multiple segmentations. 2007. 3 930
906 [16] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database
907 of human segmented natural images and its application to
908 evaluating segmentation algorithms and measuring ecolog-
909 ical statistics. In *IEEE International Conference on Computer
910 Vision*, volume 2, pages 416–423, 2001. 1 931
911 [17] S. E. Palmer. *Vision science: Photons to phenomenology*,
912 volume 1. MIT press Cambridge, MA, 1999. 3 932
913 [18] B. Peng and O. Veksler. Parameter selection for graph cut
914 based image segmentation. In *BMVC*, volume 32, pages 42–
915 44, 2008. 3 933
916 [19] J. Pont-Tuset and F. Marques. Supervised assessment of
917 segmentation hierarchies. In *Computer Vision–ECCV 2012*,
918 pages 814–827. Springer, 2012. 3 934
919 [20] X. Ren and J. Malik. Learning a classification model for seg-
920 mentation. In *Computer Vision, 2003. Proceedings. Ninth
921 IEEE International Conference on*, pages 10–17. IEEE,
922 2003. 3, 4 935
923 [21] Z. Ren and G. Shakhnarovich. Image segmentation by cas-
924 caded region agglomeration. In *Computer Vision and Pattern
925 Recognition*, 2013. 1, 2 936
926 [22] Z. Ren and G. Shakhnarovich. Image segmentation by cas-
927 caded region agglomeration. In *Computer Vision and Pat-
928 tern Recognition (CVPR), 2013 IEEE Conference on*, pages
929 2011–2018. IEEE, 2013. 6 937
930 [23] B. C. Russell, W. T. Freeman, A. A. Efros, J. Sivic, and
931 A. Zisserman. Using multiple segmentations to discover ob-
932 jects and their extent in image collections. In *Computer Vi-
933 sion and Pattern Recognition, 2006 IEEE Computer Society
934 Conference on*, volume 2, pages 1605–1614. IEEE, 2006. 3 935
935 [24] P. Salembier and L. Garrido. Binary partition tree as an effi-
936 cient representation for image processing, segmentation, and
937 information retrieval. *IEEE Transactions on Image Process-
938 ing*, 9(4):561–576, 2000. 1 939
939 [25] J. Shi and J. Malik. Normalized cuts and image segmen-
940 tation. *Pattern Analysis and Machine Intelligence, IEEE
941 Transactions on*, 22(8):888–905, 2000. 1, 6 940
942 [26] D. Varas, M. Alfaro, and F. Marques. Multiresolution hier-
943 archy co-clustering for semantic segmentation in sequences
944 with small variations. In *International Conference on Com-
945 puter Vision*, 12/2015 2015. 1 941
946 [27] C. Xu, S. Whitt, and J. J. Corso. Flattening supervoxel
947 hierarchies by the uniform entropy slice. In *Computer Vi-
948 sion (ICCV), 2013 IEEE International Conference on*, pages
949 2240–2247. IEEE, 2013. 1, 3 949
950 [28] Y. Yu, C. Fang, and Z. Liao. Piecewise flat embedding for
951 image segmentation. 2015. 6 950
952 953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971