# In the Shadows, Shape Priors Shine: Using Occlusion to Improve Multi-Region Segmentation

Yuka Kihara Ricoh Co., Ltd.\* yuuka.kihara@nts.ricoh.co.jp Matvey Soloviev Cornell University msoloviev@cs.cornell.edu Tsuhan Chen Cornell University tsuhan@cornell.edu

# Abstract

We present a new algorithm for multi-region segmentation of 2D images with objects that may partially occlude each other. Our algorithm is based on the observation that human performance on this task is based both on prior knowledge about plausible shapes and taking into account the presence of occluding objects whose shape is already known – once an occluded region is identified, the shape prior can be used to guess the shape of the missing part. We capture the former aspect using a deep learning model of shape; for the latter, we simultaneously minimize the energy of all regions and consider only unoccluded pixels for data agreement.

Existing algorithms incorporating object shape priors consider every object separately in turn and can't distinguish genuine deviation from the expected shape from parts missing due to occlusion. We show that our method significantly improves on the performance of a representative algorithm, as evaluated on both preprocessed natural and synthetic images. Furthermore, on the synthetic images, we recover the ground truth segmentation with good accuracy.

# 1. Introduction

The objective of multi-region segmentation of an image is to determine the boundaries of a number of regions containing objects of interest. Unlike the single-region case, in which the boundary unambiguously defines a partition of the image domain (and so every pixel can be associated with the object or the background), region membership becomes ambiguous when multiple regions are considered and are allowed to intersect. While several methods for multi-region segmentation have been presented in the past [24, 25, 26, 27, 31], most of these algorithms generate a true partition of the image into regions (including the background), so that every image point is assigned to a unique region. The scenario where regions may overlap arises naturally when trying to infer the true outlines of a number of objects in a two-dimensional picture of a 3D scene, which may partially occlude each other. Human observers have been known to vastly outperform known algorithms in this setting [28, 29]; folk wisdom and previous studies [1, 2, 3, 4, 5, 6] partially attribute this gap to extensive explicit (having seen many objects, familiarity with laws of physics) and implicit (neurovisual tendency towards completing straight lines, constant-curvature curves and many more complex features) prior knowledge about plausible object shapes. Yet, as the following pair of figures demonstrates, this is not the whole story:



On the right, we clearly see a figure "8", partially obscured by a black circle in front of it. The figure on the left, though, could be resolved either in the same fashion (especially if the background is not a uniform color), or as an unusually-shaped figure "3", with most observers, human or algorithmic, likely to lean towards the latter.

This suggests that besides prior knowledge about the shape of objects, knowledge of the presence and shape of occluding objects is also relevant to segmentation performance – failure of the image to manifest the expected shape of a figure "8" due to the presence of an overlapping shape counts as weaker evidence against the presence of a figure "8" than failure due to a gap or perhaps noise. In fact, it is precisely in the regions that can be established to be occluded that prior knowledge about shape attains its full significance – in the absence of data on occluded pixels, whether a certain class of object may be identified in a region depends on whether the exposed parts may be extended into an occluded region in a new way that creates a shape considered plausible for that class.

<sup>\*</sup>Work done while the first author was visiting Cornell University.

In this paper, we will explore to what extent taking both presence of other objects and prior knowledge about shapes into account may allow us to improve upon existing multi-region segmentation algorithms. To this end, we shall introduce an algorithm for multi-region segmentation that optimizes both for adherence to a machine-learned prior notion of plausible shapes (the shape prior) and agreement with the image data where it is appropriate – that is, where the object presumed to lie in the region is not occluded by some other object situated in front of it. In particular, this does entail simultaneously evolving multiple partially overlapping region boundaries. To our knowledge, no previous attempts [1, 2, 3, 4, 5, 6, 13, 14, 15, 19, 20] at multi-region segmentation attempted to employ shape priors in this fashion.

Our "plausible shapes" will be drawn from one of a number of classes of familiar objects, each represented using the Shape Boltzmann Machine deep learning model (introduced in Section 4.1). We interpret the problem as an extension of Nitzberg, Mumford and Shiota's *Segmentation with Depth* [7, 8] (Section 3), formulated with a data-driven term that favors uniformly colored regions of the input region and a shape term calculated by the machine learning model from a spectral descriptor representation of the region; we introduce a new way to combine and optimize for both of them simultaneously using the NMS functional in Section 4.2.

### 2. Related Work

As previously mentioned, the segmentation with depth problem was previously formulated by Nitzberg, Mumford and Shiota in [7, 8] using a variational formulation. The NMS model consists of a sum of a data driven term of each region, which favours segmenting curves close to the boundaries of visible parts, and a shape constraint term that constrains the possible shape contour of objects inside the image, including occluded parts. Segmentation with depth thus determines the boundaries of overlapping objects, taking into account the shape and the relative distance(depth) order of the objects, based on intensity distributions in the object regions. Numerical methods for minimizing the NMS model have been presented in [9, 10, 11]. In [12], multiple shape prior segmentation task using graph cuts has been presented. This method is capable of handling overlap by allowing for a pixel to have multiple object memberships (labels) and simultaneously segmenting multiple objects. They define the shape prior energy using a discrete version of the shape distance proposed in [11] for the level sets framework, and incorporate this energy into the graph via terminal edge weights.

However, all these works consider very limited shape priors such as contour smoothness or directly known objects, each class consist of only one image. It is not adequate to represent the shape prior encoding more complex shape variation. Recently, deep learning models [13, 14] are attractive for their well performance in modeling highdimensional richly structured data.

A deep learning model is a machine learning model that considers multiple levels of representation and abstraction that help to make sense of data such as images, sound, and text. One such model, called a Shape Boltzmann Machine (SBM) [14], is a particular form of Deep Boltzmann Machine (DBM) [13] proposed specifically for the task of modeling binary object shapes. An SBM is trained in a generative manner: given a set of shapes, the goal is to learn a probabilistic model that models the given shapes accurately and can generalize to unseen instances of the multiple shape categories. Due to its structure, SBM is successfully reduce the number of first layer parameters and facilitate efficient learning for smaller datasets while preserving its ability.

We use the SBM to incorporate information about shape in the NMS formulation. The shapes can be a set of different familiar objects, but train without information about image class. The learned model parameters implicitly define a probability distribution over all possible binary shapes in SBM; combining this prior with a variational segmentation model will necessitate a probabilistic rather than level-set representation of shape, which we will introduce in 4.2.

## 3. Overview: NMS Segmentation with Depth

We shall briefly review the formulation of segmentation with depth by Nitzberg, Mumford and Shiota[7].

Let  $\mathbf{u}: \Omega \to \mathbb{R}^1$  be a two-dimensional greyscale image defined on image domain  $\Omega \subset \mathbb{R}^2$ .

We make the following assumptions about the scene:

- 1. No object to be recognised is completely occluded.
- Every object is of approximately uniform brightness, and any two objects in the scene have different brightness.
- 3. The objects are not interlaced, so in particular there is a well-defined ordering by distance from the observer.

Note that by ignoring invisible objects, all three assumptions are easily satisfied by considering depth images of spatially well-separated objects.

Let  $R_1, \ldots, R_n$  be the regions occupied by objects, ordered by increasing distance from the observer (i.e. depth). Each region  $R_i$  may be partially occluded by the regions preceding it in the ordering (but not completely so, by assumption 1). We shall denote the visible parts of the above objects as  $R'_1, \ldots, R'_n$  respectively; they may be formally defined as

$$\begin{array}{rcl} R_1' &=& R_1 \\ R_i' &=& R_i - \bigcup_{j < i} R_j \text{ for } i > 1. \end{array}$$

We may denote the background as  $R'_{n+1} = \Omega - \bigcup_{j \le n} R_j$  for the sake of consistency.

**Definition 3.1.** (*NMS*) Given an image  $f : \Omega \to \mathbb{R}$  and a number n of recognisable objects, the objective of segmentation with depth is to determine

- an ordering O<sub>1</sub> < ... < O<sub>n</sub> of recognisable objects in the image by distance from the viewer;
- the average greyscale intensity c<sub>i</sub> ∈ ℝ of the region occupied by each object O<sub>i</sub> and
- the shape R<sub>i</sub> ⊆ Ω of the region occupied by each object O<sub>i</sub>.

The *NMS functional* provides an objective function for comparing candidate solutions for the third point, given solutions for the first two. In its general form, it may be written as follows:

$$E_{\text{NMS}} = \sum_{i=1}^{n+1} \int_{R'_i} |\mathbf{u}(x,y) - c_i|^2 dx dy + E_{\text{shape}}.$$
 (1)

The  $E_{\text{shape}}$  term of  $E_{NMS}$  serves to penalise implausible shapes of regions. In [7], the authors set

$$E_{\text{shape}} = \sum_{i=1}^{N} \oint_{\partial R_i \cap \Omega} [\alpha + \beta \psi(\kappa)] ds$$

As this contour integral is taken along the boundary of the original unoccluded shape, evaluating the function inevitably entails reconstructing invisible parts of regions' boundaries. Depending on the setting of the parameters  $\alpha$ and  $\beta$ , one may obtain curvature functions such as Euler's Elastica.

In our work, instead of a curvature-based boundary energy function, we use the evaluation of a deep learning model trained on global and local features of "good"/"known" shapes, which we shall briefly introduce in the next section, plus a denoising term based on the total variance norm as in [15]. The top-level segmentation algorithm then takes the following form:

- 1. By running k-means on the histogram of the input depth image u, determine the number and average intensity of objects present. As intensity represents depth in the image, this also gives a canonical depth ordering.
- 2. For each object i, find a small uniform region of the average colour determined for it initialise  $R_i$  with it.
- 3. Perform gradient descent on a relaxed version of the NMS functional using the Split Bregman method (section 4.3) to find a local minimum for  $\{R_i\}_{1 \le i \le n}$ .

# 4. Algorithm

#### 4.1. The Shape Boltzmann Machine

In our work, we implement the shape energy term using a *Shape Boltzmann machine* (SBM). (For an in-depth introduction of SBMs and DBMs, we refer the reader to [14] and [13] respectively.) Our SBM employs two layers of hidden variables. We shall denote by v the vector of visible units that represents the (input) binary shape image, and by  $h^1$ and  $h^2$  be the lower and higher binary hidden units respectively.

The first layer  $h^1$  consists of several disjoint subsets  $h_1^k$  of same size m. Each of these hidden units has a restricted receptive field and only connects to a subset of the visible units, but sharing weights between the sets of hidden units and visible units  $W^1 = W_k^1$ ,  $k \in K$ . The visible units are local square patches having the same size n and overlapping its neighbor by d pixels along the boundaries.

Similar to DBMs, there are no within-layer connections. The (output) energy of the state  $\{v, h^1, h^2\}$  is then defined as follows:

$$E_{SBM}(\mathbf{v}, \mathbf{h}^{1}, \mathbf{h}^{2}; \theta)$$

$$= -\sum_{k \in K} \mathbf{v}_{\mathbf{k}}^{\mathrm{T}} \mathbf{W}^{1} \mathbf{h}_{\mathbf{k}}^{1}$$

$$-\mathbf{b}^{\mathrm{T}} \mathbf{v} - \mathbf{c}^{1\mathrm{T}} \mathbf{h}^{1} - \mathbf{h}^{1\mathrm{T}} \mathbf{W}^{2} \mathbf{h}^{2} - \mathbf{c}^{2\mathrm{T}} \mathbf{h}^{2} \quad (2)$$

Here,  $\theta = { \tilde{\mathbf{W}}^1, \mathbf{W}^2, \mathbf{c}^1, \mathbf{c}^2, \mathbf{b} }$  are the model parameters.  $\tilde{W}^1$  is defined as follows. The first term can be rewritten in the same form of DBM by some matrix manipulation:

$$\sum_{k \in K} \mathbf{v}_k^{\mathrm{T}} \mathbf{W}^1 \mathbf{h}_k^1 = \mathbf{v}^{\mathrm{T}} \tilde{\mathbf{W}}^1 \mathbf{h}^1$$
(3)

The distribution over v is given by marginalizing over the hidden variables.

$$p(\mathbf{v};\theta) = \frac{1}{Z(\theta)} \sum_{\mathbf{h}^1, \mathbf{h}^2} exp(-E(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2))$$
(4)

Here,  $Z(\theta)$  is the partition function. Although exact inference is no longer tractable, an efficient approximate learning can be carried out by using a mean-field procedure based on the property that conditional distributions  $p(\mathbf{v}|\mathbf{h}^1)$ ,  $p(\mathbf{h}^1|\mathbf{v}, \mathbf{h}^2)$ , and  $p(\mathbf{h}^2|\mathbf{h}^1)$  remain independent, as detailed in [13]. Each conditional distribution is given by sigmoid function  $\sigma(y) = 1/(1 + exp(-y))$  as follows.

$$p(h_j^1 = 1 | \mathbf{v}, \mathbf{h}^2) = \sigma(\sum_i \tilde{W}_{ij}^1 v_i + \sum_k W_{jk}^2 h_k^2 + c_j^1)$$
(5)

$$p(h_k^2 = 1 | \mathbf{h}^1) = \sigma(\sum_j W_{jk}^2 h_j^1 + c_k^2)$$
(6)

$$p(v_i = 1 | \mathbf{h^1}) = \sigma(\sum_j \tilde{W}_{ij}^1 h_j^1 + b_i)$$
(7)

Given a set of aligned binary shape images, we learn the model of SBM for shape prior segmentation which maximize likelihood  $p(\mathbf{v}; \theta)$  with respect to  $\theta$ . As proposed in [13], the model parameters  $\theta$  can be efficiently pre-trained at each layer greedily layer-by-layer. Joint training is then carried out to fine-tune the parameters and separate learning of local and global shape properties into the two hidden layers.

#### 4.2. Incorporating the Shape Prior

If we were to optimize for only the first term of our objective function 1, we would be solving an instance of a normal multi-region segmentation problem in which we partition the image into disjoint regions  $R'_i$  associated with each object – the possibly overlapping inferred boundaries  $R_i$  would be playing no role. Since parametric representations of boundaries are typically insufficiently expressive and/or badly behaved for optimization purposes, in the literature, minimization of objective functions of this form is typically handled using the level set method [27], where the interior of each region is represented as the positive locus of some continuous function of the image domain and a time parameter. For an overview of the level set method in the context of image processing, see [27].

In the level set framework, region shape is commonly represented using a *signed distance function*, that is, a function that gives the shortest distance to the nearest point on the boundary inside the region and its negative outside. The SBM used for the shape energy term, however, instead operates on a different continuous relaxation of region membership, namely a function  $\mathbf{q} : \Omega \rightarrow [0,1]$  that assigns a probability that each pixel of the image is contained within the region in question. To be able to incorporate the shape energy term, we need to find a common representation, and reformulating the data term in terms of  $\mathbf{q}$  is comparatively straightforward.

When working with such a function, a natural approach is to optimize on it directly and pick some cutoff probability if a binary notion of region membership is again required. This was first considered by Cremers, Schmidt and Barthel in [16] following [15] with the purpose of convexifying the space that optimization is performed over.

Recall (1) that for a single region and its complement, the data term of the NMS functional took the following form:

$$E_{\text{data}} = \int_{R} |\mathbf{u} - c_1|^2 dx dy + \int_{R^c} |\mathbf{u} - c_2|^2 dx dy$$

Here,  $c_1 = c_1(R)$  and  $c_2 = c_2(R)$  are the average intensities in R and  $R^c$  respectively, and so in particular depend on R.

Suggestively, we may reformulate this expression with a single integral taken over the whole image domain using the

characteristic function  $\chi_R$  of the region R as follows:

$$E_{\text{data}} = \int_{\Omega} |\mathbf{u} - c_1|^2 \chi_R + |\mathbf{u} - c_2|^2 (1 - \chi_R) dx dy.$$

As  $\mathbf{q}$  can be seen as a relaxation of  $\chi_R$ , this immediately suggests the following relaxed form in terms of  $\mathbf{q}$ :

$$E_{\text{data}}(\mathbf{q}) = \int_{\Omega} |\mathbf{u} - c_1|^2 \mathbf{q} + |\mathbf{u} - c_2|^2 (1 - \mathbf{q}) dx dy.$$
(8)

#### **4.3. Decomposed Objective Function**

We will now show how to integrate the above shape representation with the shape prior model 11 to obtain an objective function for the problem that we may perform optimization on.

Noting how the data term for one region and the background in equation 8 consist of a part associated with the region and one associated with its complement, we can decompose the data-based energy term to obtain a collection of coupled energy terms for each of the curve boundaries individually following the approach taken in [24]. We modify the NMS representation to keep track of the probability that each point is part of a given region, instead of the usual binary notion of region membership. Letting  $\mathbf{q}_i : \Omega \rightarrow [0, 1]$ be the membership probability function of the region  $R_i$ , we set

$$E_{data}(\mathbf{q}_i) \tag{9}$$

$$= \int_{\Omega} \prod_{j=1}^{i-1} (1 - \mathbf{q}_j) \left( |\mathbf{u} - c_i|^2 \mathbf{q}_i + (1 - \mathbf{q}_i) \Phi_i \right) dx dy,$$

where  $\Phi_i$  is

$$\Phi_i = \sum_{j=i+1}^{n+1} \left( |\mathbf{u} - c_j|^2 \mathbf{q}_j \prod_{k=2}^{j-1} (1 - \mathbf{q}_k) \right) \quad (10)$$

with  $\mathbf{q}_{n+1}$  understood to be 1 everywhere; the product of  $\mathbf{q}_i$ s and  $(1-\mathbf{q}_i)$ s can be seen to be the natural generalisation of the characteristic function of the visible region  $R'_i$  to the probabilistic relaxation.

We combine the shape prior term with this data term and reformulate our objective function (1) as follows:

$$E_{SPSD} = \sum_{i=1}^{N} \left( E_{data}(\mathbf{q}_i) + \mu E_{shape}(\mathbf{q}_i) \right)$$
$$= \sum_{i=1}^{N} \int_{\Omega} |\mathbf{u} - c_i|^2 \mathbf{q}_i + \Phi_i (1 - \mathbf{q}_i) dx dy$$
$$- \mu (\mathbf{q}_i^{\mathrm{T}} W^1 \mathbf{h}^1 - \mathbf{q}_i^{\mathrm{T}} \mathbf{b} - \mathbf{c}^{1\mathrm{T}} \mathbf{h}^1 - \mathbf{h}^{1\mathrm{T}} \mathbf{W}^2 \mathbf{h}^2 - \mathbf{c}^{2\mathrm{T}} \mathbf{h}^2)$$
$$+ \nu ||\nabla \mathbf{q}||_e. \tag{11}$$

Note that we consider every region simultaneously while performing gradient descent. In particular, at any step, only the part of a given region that is taken to be unoccluded is compared to its mean luminance for the data agreement term; e.g. a hypothetical completely occluded region would then simply converge to the nearest plausible shape according to the shape prior with no regard for data agreement.

#### 4.4. Minimizing the Energy Function

When the learning model parameters are known,  $E_{SPSD}$ has two kind of unknowns: the shapes  $q_i$  and the SBM related hidden units  $h^1$  and  $h^2$ . Instead of addressing both together, we use an alternating minimization procedure. Each layer of hidden units can be computed by mean-field approximate inference, just as done for DBM [15]. Observe that energy functional 11 with respect to each shape  $q_i$  is a convex functional, and hence can be solved using the Split Bregman method to obtain a global minimizer. We apply the Split Bregman method [17] to solve the minimization problem of step 5 in Algorithm 1, adapting the implementation in [18] with respect to q, and perform exhaustive search to find the minimizing translation and scaling of W.

Algorithm 1: Multi-region segmentation with depth
<b>Input:</b> an image <i>u</i> , model parameters
$\theta = { {\tilde{\mathbf{W}}^1, \mathbf{W}^2, \mathbf{c}^1, \mathbf{c}^2, \mathbf{b} } $ , and a test image <b>u</b>
1 Initialize $\{\mathbf{q}_i^0\}_{i=1}^{N-1}$
repeat
$2     k \leftarrow k+1$
for $i = 1 \dots N - 1$ do
3 $\mathbf{h}^{1} \leftarrow \sigma((\mathbf{q}_{i}^{k-1})^{\mathrm{T}} \tilde{W}^{1} + W^{2} h^{2} + c^{1})$
4 $\mathbf{h}^2 \leftarrow \sigma(\mathbf{h}^{1T}\mathbf{W}^2 + \mathbf{c}^2)$
5 $\mathbf{q}_i^k \leftarrow$
$\arg\min_{\mathbf{q},\mathbf{W}}(E_{data}(\mathbf{q}) - \mathbf{q}^{\mathrm{T}}\tilde{\mathbf{W}}^{1}h^{1} - \mathbf{q}^{\mathrm{T}}b)$
until $\sum_{i=1}^{N-1} \ \mathbf{q}_i^{\ k} - \mathbf{q}_i^{\ k-1}\  < \epsilon;$
7 return $\{q_i\}_{i=1}^{N-1}$

The initialization step in Algorithm 1 is performed by running one iteration of the Split Bregman minimization without a shape prior, i.e. setting  $E_{SBM}(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2; \theta) = 0$ . We have a  $O(n^2)$  dependency on number of regions, as the shape prior is applied independently, and each  $E_{data}(\mathbf{q}_i)$  in equation (9) takes O(n) to calculate naively. This may easily be improved to O(n) by carrying through partial products. The dependency on image size is also  $O(n^2)$ , as both the Split Bregman iteration and our updates of  $\mathbf{q}$  are local.

## 5. Experiments

All experiments were run in MATLAB on a PC with a 2.30 GHz Intel Core i7-3610QM Processor and 8GB of RAM. To implement the Split Bregman arg min in Algorithm 1, we used the MATLAB native code wrapper to interface with a modified version of the code provided by T. Goldstein as a supplement to [18].

We evaluate the effectiveness of our algorithm on several datasets which consist of binary segmentation masks representing an object silhouette (Fig. 1(top)). Through the evaluation, those binary images are cropped and normalized to the specific size, and about half of the images are used for training, and the rest of them for testing. In the SBM training phase, pre-training requires 3000 and 1000 epochs for the first and second layers. In addition, global training is performed for 1000 epochs.

As corresponding algorithms using Split Bregman optimization with Markov Random Fields (MRF), Factor Analysis [30], and the Restricted Boltzmann Machine (RBM) have been investigated in [14], we restrict ourselves to evaluating our algorithm's performance on each input in comparison to that of the single-region shape prior segmentation algorithm of Chen et al. (CVPR 2013) [15], with the DBM model it employs replaced with our SBM model to ensure that only the effect of treating occlusion differently is captured. We employ this single-region algorithm for multiregion segmentation in a natural way by considering every object separately in turn and treating the rest of the image as the "background". In the tables, this algorithm is referred to as "SP (single)". We also consider the results obtained by an analogous algorithm that does not employ the machinelearned shape priors at all (referred to as "no SP") for perspective.

We evaluate the algorithms on both synthetic images and real images. The synthetic images are generated by overlapping two or three binary shapes selected uniformly at random from the test set. We also added Gaussian noise to each image. Fig. 1 shows some examples of the synthetic images used in the experiment.

For the synthetic images, the accuracy of the segmentation can be evaluated by comparing it to the ground truth segmentation, i.e. the individual silhouettes that were composed to form the test image. Since each segmented region corresponding to  $\mathbf{q_i}$  has a pixel value in [0,1], where each value represents the probability that the pixel is in the region, we simply take as the region those pixels having value larger than 0.5 to compare with the binary segmentation mask. We use two common metrics to assess segmentation quality: the average pixel accuracy (AP), of foreground and background classification, and the foreground intersectionover-union score (IoU), defined as  $O_{iou}(S_p, S_q) = \frac{S_p \cap S_p}{S_p \cup S_p}$ where  $S_p$  and  $S_q$  are the sets of ground truth and predicted foreground pixels.

#### 5.1. Single-Class Dataset

We evaluate the basic performances on two datasets.

**Caltech-101 motorbikes**: The first dataset we investigated was 798 motorbike silhouettes from Caltech-101 dataset [22], capturing motorbike from side. We used 399 images



Figure 1. Some examples of shape dataset (top) and the synthetic images (bottom). The synthetic image is generated by locating two or three resized binary shapes and adding noise. The combinations of the shapes used in the synthetic images are randomly selected.

for training, and the rest of them for testing. The images are cropped and normalized to  $64 \times 64$  pixels. We trained an SBM with d = 4, and 1200 and 50 units for  $h^1$  and  $h^2$  respectively.

**Weizmann Horses:** The second dataset is the Weizmann horse dataset [23] consisting of 328 images, all of horses facing left, but with a high variety of poses. We created a training set of 164 images from this dataset, cropped and normalized to  $32 \times 32$  pixels. We trained an SBM with d = 4, and 2000 and 100 units for  $h^1$  and  $h^2$  respectively.

Example images of results for motorbikes and Weizmann horses are shown in Fig 2. The shape prior based algorithms can be seen to benefit from the power of the SBM model: even if one of the main components of the object, such as a handle of the motorbike or a leg of the horse, is missing, the proposed method still determine an proper boundary of the object recovering its plausible shape. Red indicates the pixels within the first object from the front, green indicates the second object from the front, and blue indicates the third object from the front.

The qualitative and quantitative comparison of our algorithm with "SP (single)" and "no SP" can be seen in Fig. 3 and Table 1. Testing was performed with 399 motorbike images and 164 horse images.

We expected the single-region approach to suffer in occluded regions as  $E_{data}$  is raised by the mismatch in pixel intensity. The measurements appear to confirm this prediction – our method exhibits significantly improved scores for background regions across the board, and visual inspection confirms that this is largely due to improved accuracy in occluded parts. The keen observer may note that using a shape prior results in slight loss of accuracy for the front-most region. This is unsurprising: without the shape energy term, the contour is free to evolve to exactly match the uniformly colored region, but artifacts of the SBM shape model may lead to slight deviations having lower total energy.



Figure 2. Segmentation results with a single object class for two (a-d) and three (e-h) regions. Some input images, together with the respective ground truth shapes, are shown in (a)(c) and (e)(g). The composite segmentation outputs (with the regions being rendered, from front to back, in red, green and blue), along with the recovered shapes, can be seen in (b)(d) and (f)(h).



Figure 3. Segmentation comparison of our method (multi-region segmentation) vs. single-region segmentation running several times: Original synthetic images are shown in (a),(d). (b)(e) and (c)(f) show the segmentation results and the corresponding shapes of single-region segmentation, our method, respectively. Our results are encircled by red line.

## 5.2. Multi-Class dataset

Next, we shows the results of multi region segmentation with multiclass shape prior. We tested how well our method infers highly occluded object shape without any information about the category of each object, leaving it to the shape model to discover the most likely object categories from input image. For this purpose, we trained SBM on a combination of the images of motorbike, revolver, and car of side view from Caltech-101 dataset [22]. The training dataset contains 60 cars in profile, 40 revolvers, and 300 motor bikes for a total 400 images. The images are cropped and normalized to  $64 \times 32$  pixels. A SBM with d = 4, and

method	AP		IoU		AP			IoU		
	region1	region2	region1	region2	region1	region2	region3	region1	region2	region3
no SP	100	88.67	100	78.76	100	92.77	71.48	99.99	86.23	49.51
SP (single)	98.94	89.35	98.27	80.95	98.79	92.88	82.87	98.08	87.24	71.39
SP (multi)	98.86	94.04	98.13	89.52	97.88	96.25	86.98	96.64	93.24	78.41
no SP	99.92	80.95	99.83	65.70	99.97	81.98	67.46	99.94	69.27	41.94
SP (single)	98.69	82.58	97.61	69.12	98.29	82.70	69.97	96.88	70.42	47.91
SP (multi)	98.36	87.62	97.01	77.96	97.65	86.00	71.39	95.71	75.89	50.36

Table 1. Results on the single class dataset: bike(top), horse(bottom). All standard deviations are bounded above by  $\pm 0.23$ .

method	AP		IoU		AP			IoU		
	region1	region2	region1	region2	region1	region2	region3	region1	region2	region3
no SP	100	84.35	100	71.89	99.92	85.30	63.10	99.91	73.80	32.53
SP (single)	98.56	85.88	97.93	76.74	98.43	85.59	70.58	97.59	76.96	57.96
SP (multi)	98.22	86.39	98.22	86.39	98.03	92.04	77.96	97.00	86.76	66.35

Table 2. Results on the multiclass dataset with three object categories from Caltech-101. All std. deviations are bounded above by  $\pm 0.16$ .



Figure 4. Segmentation results for multiple object classes. Input images and ground truth segmentations shown in (a),(c), (e) and (g). Segmentation results are in (b),(d), (f) and (h). (h) depicts an instance of segmentation failing: the third object, a partially occluded motorbike, is recovered as a revolver.

dataset	sin	gle	multi			
	2 regions	3 regions	2 regions	3 regions		
bike	$2.39(\pm 0.80)$	$5.43(\pm 1.00)$	$3.50(\pm 3.04)$	$7.63(\pm 3.32)$		
horse	$1.74(\pm 0.62)$	$4.03(\pm 1.74)$	$1.05(\pm 0.73)$	$4.36(\pm 2.08)$		
Caltech	$2.28(\pm 0.81)$	5.31(±0.44)	$1.45(\pm 1.29)$	$8.32(\pm 3.37)$		

Figure 5. Average runtime per instance for the three synthetic datasets, measured on our test system, in seconds.

1500 and 300 units for  $h_1$  and  $h_2$  was jointly trained without information about image class.

Example segmentation results can be seen in Fig. 4; quantitative results from 599 synthetic images can be found

in Table 2. We observe that our algorithm largely identifies the right object category and converges to the appropriate shape. There are, however, some cases where the occlusion of the shape is too heavy to determine the plausible object category. In the false example in Fig. 4, a partially occluded motorbike is recovered as a revolver.

#### 5.3. Enhanced Real-World Images

Next, we tested the proposed method on real images. In application, the target object often has similar shapes in different scales. The data set we tested here contains partially occluded cars with different direction, scale, and location. We used 60 side car shapes (out of a total 123) from the Caltech 101 dataset, with all cars facing left but not necessarily being of the same type. In order to segment a car facing the other direction simultaneously, we construct another 60 training shapes by flipping the previous training shapes horizontally. We cropped and normalised all shapes to  $64 \times 32$  pixels, and trained an SBM with 1200 units in the first and 50 units in the second hidden layer on this extended training set with a total 120 shapes. The initial segmentation of visible part of the object is manually given. Since the aspect ratio of the car shapes is changing in the training phase, we normalized input images so that the size of the car located forward will be ratio of 64:32 in segmentation process. Here, too, we find that our algorithm outperforms the reference in the same way (Fig. 6).

We have found that the performance depends on the occlusion ratio. If more than half of the shape is occluded, our algorithm often mistakes the the visible part for one belonging to a smaller copy of the object. Some consecutive frames of a video of a moving car have been segmented in Fig. 7 to show how segmentation performance declines as more and more of an object is occluded.



Figure 6. Qualitative test results on enhanced real-world images. Top row: Input. Middle row: Results of our algorithm, overlaid onto the original image. Bottom row: Results of the single-region shape prior algorithm.



Figure 7. Qualitative test results on a few subsequent frames of enhanced real-world images of a moving car.

One possible approach to overcome this limitation is to integrate *context*, such as scale (objects further away should be smaller), orientation (should be similar each other among existing objects), or physical plausibility, in our algorithm.

## 6. Conclusion

We have presented a variational algorithm for multiregion segmentation of partially overlapping objects. Our algorithm uses prior information about object shapes modeled by a Shape Boltzmann Machine (SBM), and integrates it with a data-driven term. The data-driven term is made to depend only on image pixels assumed to be relevant, based on a probabilistic representation of region membership compatible with the SBM. On one hand, this enables it to deal with noise and damage in visible parts in the usual way; on the other hand, it also freely infers the shape of parts that are missing in the image due to being occluded, without having to compete with irrelevant data originating from occluding objects. The latter enables it to exhibit significant improvements over a variational algorithm representative of the state of the art, which treats occluding objects as damage.

We believe that our data supports the thesis that accounting for occluding objects is crucial to approaching human performance at multi-region segmentation. Since many subpar segmentations produced by our algorithm exhibit shapes that get good scores from the SBM but do not actually depict plausible silhouettes of the object in question, performance may likely be improved further by employing a more advanced model for the shape prior.

## Acknowledgements

The authors were partially supported by the following grants: NSF CCF-1214844, MURI FA9550-12-1-0040 and ARO W911NF-09-1-0281. We would like to thank the Cornell AMP Lab, Kevin Matzen, Kyle Wilson, the Cornell Graphics and Vision group and the anonymous reviewers for support and helpful comments.

# References

- M. E. Leventon, W. E. L. Grimson, and O. Faugeras. Statistical shape influence in geodesic active contours. In Proc. Conf. Comput. Vis. Pattern Recognit., Vol. 1. 2000, pp. 316323.
- [2] M. Rousson and N. Paragios. Shape Priors for Level Set Representations. In Proc. Eur. Conf. Comput. Vis., 2002, pp. 7892.
- [3] A. Tsai, A. Yezzi, Jr., W. Wells, C. Tempany, D. Tucker, A. Fan, W. E. Grimson, and A. Willsky. A shape-based approach to the segmentation of medical imagery using level sets. IEEE Trans. Med. Imag., Vol. 22, no. 2, pp. 137154, Feb. 2003.
- [4] M. Rousson and D. Cremers. Efficient Kernel Density Estimation Of Shape And Intensity Priors For Level Set Segmentation. In Proc. MICCAI, Vol. 3750. 2005, pp. 757764.
- [5] D. Cremers, S. J. Osher, and S. Soatto. Kernel Density Estimation and Intrinsic Alignment for Shape Priors in Level Set Segmentation. Int. J. Comput. Vis., Vol. 69, no. 3, pp. 335351, 2006.
- [6] D. Cremers. Nonlinear Dynamical Shape Priors for Level Set Segmentation. J. Sci. Comput., Vol. 35, nos. 23, pp. 132143, 2008.
- [7] M. Nitzberg and D. Mumford. The 2.1-d Sketch. In ICCV, 1990.
- [8] M. Nitzberg, D. Mumford, and T. Shiota. Filtering, segmentation and depth. Lecture Notes in Computer Science, 662, 1993.
- [9] W. Zhu, T. Chan, and S. Esedoğlu. Segmentation with depth: a level set approach. SIAM J. Sci. Comput., 28 (2006), pp. 19571973.
- [10] S. Esedoğlu and R. March. Segmentation with depth but without detecting junctions. J. Mathematical Imaging and Vision, 18(1):715, 2003.
- [11] Sheshadri R. Thiruvenkadam, Tony F. Chan, and Byung-Woo Hong. Segmentation under Occlusions Using Selective Shape Prior. Scale Space and Variational Methods in Computer Vision. Springer Berlin Heidelberg, 2007. 191-202.
- [12] Nhat Vu and B. S. Manjunath. Shape prior segmentation of multiple objects with graph cuts. In CVPR, 2008.
- [13] R. Salakhutdinov and G. Hinton. Deep Boltzmann machines. In AISTATS, 2009.

- [14] S. M. A. Eslami, N. Heess, and J. Winn. The Shape Boltzmann machine: a Strong Model of Object Shape. In CVPR, 2012.
- [15] F. Chen, H. Yu, R. Hu, and X. Zeng. Deep learning shape priors for object segmentation. In CVPR, 2013.
- [16] D. Cremers, F. Schmidt, and F. Barthel, Shape priors in variational image segmentation: convexity, lipschitz continuity and globally optimal solutions, in CVPR, 2008.
- [17] T. Goldstein and S. Osher. The Split Bregman Method for L1-Regularized Problems. SIAM J. Image Sciences, 2(2):323343, 2009.
- [18] T. Goldstein, X. Bresson, and S. Osher. Geometric applications of the split Bregman method: segmentation and surface reconstruction, J Sci Comput, 45: 272293, 2010
- [19] Victor Adrian Prisacariu, Aleksandr V. Segal, and Ian Reid. Simultaneous monocular 2D segmentation, 3D pose recovery and 3D reconstruction. ACCV 2012. Springer Berlin Heidelberg, 2013. 593-606.
- [20] Zhi Yang, Yu Kong, and Yun Fu. Decomposed contour prior for shape recognition. Pattern Recognition (ICPR), 2012 21st International Conference on. IEEE, 2012.
- [21] Y. Freund and D. Haussler. Unsupervised Learning of Distributions on Binary Vectors Using Two Layer Networks. Advances in Neural Information Processing Systems 4, 1992.
- [22] Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 28, Issue 4 (2006): 594-611.
- [23] E. Borenstein and S. Ullman. Class-Specific, Top-Down Segmentation. ECCV 2002. Springer Berlin Heidelberg, 2002. 109-122.
- [24] Abdol-Reza Mansouri, Amar Mitiche, and Carlos Vzquez. Multiregion competition: A level set extension of region competition to multiple region image partitioning. Computer Vision and Image Understanding 101.3 (2006): 137-150.
- [25] Carlos Vazquez, Amar Mitiche, and Robert Laganiere. Joint multiregion segmentation and parametric estimation of image motion by basis function representation and level set evolution. Pattern Analysis and Machine Intelligence, IEEE Transactions on 28.5 (2006): 782-793.

- [26] Ismail Ben Ayed, Amar Mitiche, and Ziad Belhadj. Polarimetric image segmentation via maximumlikelihood approximation and efficient multiphase level-sets. Pattern Analysis and Machine Intelligence, IEEE Transactions on 28.9 (2006): 1493-1500.
- [27] Amar Mitiche and Ismail Ben Ayed. Variational and level set methods in image segmentation. Vol. 5. Springer Science and Business Media, 2010.
- [28] R.A. Rensink and J. T. Enns. Early completion of occluded objects. Vision Research 38 (1998): 2489-2505.
- [29] J. Drivera et al. Segmentation, attention and phenomenal visual objects. Cognition 80.61 (2001): 95.
- [30] S. M. A. Eslami and C. K. I. Williams. Factored Shapes and Appearances for Parts-based Object Understanding. In BMVC, pages 18.118.12, 2011.
- [31] Qi Song et al. Simultaneous Searching of Globally Optimal Interacting Surfaces with Shape Priors. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2010. IEEE, 2010.