Affinity CNN: Learning Pixel-Centric Pairwise Relations for Figure/Ground Embedding

Michael Maire TTI Chicago

mmaire@ttic.edu

Takuya Narihira Sony Corp. takuya.narihira@jp.sony.com Stella X. Yu UC Berkeley / ICSI stellayu@berkeley.edu

Abstract

Spectral embedding provides a framework for solving perceptual organization problems, including image segmentation and figure/ground organization. From an affinity matrix describing pairwise relationships between pixels, it clusters pixels into regions, and, using a complex-valued extension, orders pixels according to layer. We train a convolutional neural network (CNN) to directly predict the pairwise relationships that define this affinity matrix. Spectral embedding then resolves these predictions into a globallyconsistent segmentation and figure/ground organization of the scene. Experiments demonstrate significant benefit to this direct coupling compared to prior works which use explicit intermediate stages, such as edge detection, on the pathway from image to affinities. Our results suggest spectral embedding as a powerful alternative to the conditional random field (CRF)-based globalization schemes typically coupled to deep neural networks.

1. Introduction

Systems for perceptual organization of scenes are commonly architected around a pipeline of intermediate stages. For example, image segmentation follows from edge detection [12, 1, 2, 7, 4]; figure/ground, occlusion, or depth layering follows from reasoning over discrete contours or regions [27, 16, 21, 36, 18] with some systems also reliant on motion cues [30, 15, 32, 31]. This trend holds even in light of rapid advancements from designs centered on convolutional neural networks (CNNs). Rather than directly focus on image segmentation, recent CNN architectures [14, 3, 28, 4] target edge detection. Turaga *et al.* [33] make the connection between affinity learning and segmentation, yet restrict affinities to be precisely local edge strengths. Pure CNN approaches for depth from a single image do focus on directly constructing the desired output [9, 8]. However, these works do not address the problem of perceptual grouping without fixed semantic classes.

We engineer a system for simultaneous segmentation and figure/ground organization by directly connecting a CNN to



Figure 1. System architecture. We send an image through a CNN which is trained to predict the grouping and ordering relations between each of the *n* pixels and its neighbors at *k* displacements laid out in a fixed stencil pattern. We assemble these $n \times 2k$ pixel-centric relations into a sparse $n \times n$ complex affinity matrix between pixels, each row indicating a pixel's affinity with others. Shown above is the row for the pixel at the center of a log-polar sampling pattern; its positive/negative relations with neighbors are marked by red/cyan squares overlaid on the image. We feed the pairwise affinity matrix into Angular Embedding for global integration, producing an eigenvector representation that reveals figure-ground organization: we know not only which pixels go together, but also which pixels go in front.

an inference algorithm which produces a globally consistent scene interpretation. Training the CNN with a target appropriate for the inference procedure eliminates the need for hand-designed intermediate stages such as edge detection. Our strategy parallels recent work connecting CNNs and conditional random fields (CRFs) for semantic segmentation [6, 20, 35]. A crucial difference, however, is that we handle the generic, or class independent, image partitioning problem. In this context, spectral embedding, and specifically Angular Embedding (AE) [37, 38], is a more natural inference algorithm. Figure 1 illustrates our architecture.

Angular Embedding, an extension of the spectral relaxation of Normalized Cuts [29] to complex-valued affinities, provides a mathematical framework for solving joint group-

1

ing and ranking problems. Previous works established this framework as a basis for segmentation and figure/ground organization [22] as well as object-part grouping and segmentation [24]. We follow the spirit of [22], but employ major changes to achieve high-quality figure/ground results:

- We reformulate segmentation and figure/ground layering in terms of an energy model with pairwise forces between pixels. Pixels either bind together (group) or differentially repel (layer separation), with strength of interaction modulated by confidence in the prediction.
- We train a CNN to directly predict all data-dependent terms in the model.
- We predict interactions across multiple distance scales and use an efficient solver [23] for spectral embedding.

Our new energy model replaces the ad-hoc boundarycentric interactions employed by [22]. Our CNN replaces hand-designed features. Together they facilitate learning of pairwise interactions across a regular stencil pattern. Choosing a sparse stencil pattern, yet including both shortand long-range connections, allows us to incorporate multiscale cues while remaining computationally efficient.

Section 2 develops our model for segmentation and figure/ground while providing the necessary background on Angular Embedding. Section 3 details the structure of our CNN for predicting pairwise interaction terms in the model.

As our model is fully learned, it could be trained according to different notions of segmentation and figure/ground. For example, consistent definitions for figure/ground include true depth ordering as in [9], object class-specific foreground/background separation as in [24], and boundary ownership or occlusion as in [27, 13, 22]. We focus on the latter and define segmentation as a region partition and figure/ground as an ordering of regions by occlusion layer. The Berkeley segmentation dataset (BSDS) provides ground-truth annotation of this form [25, 13]. We demonstrate segmentation results competitive with the state-ofthe-art on the BSDS benchmark [11], while simultaneously generating high-quality figure/ground output.

The occlusion layering interpretation of figure/ground is the one most likely to be portable across datasets; it corresponds to a mid-level perceptual task. We find this to be precisely the case for our learned model. Trained on BSDS, it generates quite reasonable output when tested on other image sources, including the PASCAL VOC dataset [10]. We believe this to be a significant advance in fully automatic perceptual organization. Section 4 presents experimental results across all datasets, while Section 5 concludes.

2. Spectral Embedding & Generalized Affinity

We abstract the figure/ground problem to that of assigning each pixel p a rank $\theta(p)$, such that $\theta(\cdot)$ orders pixels by



Figure 2. Angular Embedding [38]. Given (C, Θ) capturing pairwise relationships between nodes, the Angular Embedding task is to map those nodes onto the unit semicircle, such that their resulting absolute positions respect confidence-weighted relative pairwise ordering (Equation 1). Relative ordering is identified with rotation in the complex plane. For node p, $\theta(p) = \arg(z(p))$ recovers its global rank order from its embedding z(p).

occlusion layer. Assume we are given estimates of the relative order $\Theta(p,q)$ between many pairs of pixels p and q. The task is then to find $\theta(\cdot)$ that agrees as best as possible with these pairwise estimates. Angular Embedding [38] addresses this optimization problem by minimizing error ε :

$$\varepsilon = \sum_{p} \frac{\sum_{q} C(p,q)}{\sum_{p,q} C(p,q)} \cdot |z(p) - \tilde{z}(p)|^2 \tag{1}$$

where C(p,q) accounts for possibly differing confidences in the pairwise estimates and $\theta(p)$ is replaced by $z(p) = e^{i\theta(p)}$. As Figure 2 shows, this mathematical convenience permits interpretation of $z(\cdot)$ as an embedding into the complex plane, with desired ordering $\theta(\cdot)$ corresponding to absolute angle. $\tilde{z}(p)$ is defined as the consensus embedding location for p according to its neighbors and Θ :

$$\tilde{z}(p) = \sum_{q} \tilde{C}(p,q) \cdot e^{i\Theta(p,q)} \cdot z(q)$$
(2)

$$\tilde{C}(p,q) = \frac{C(p,q)}{\sum_{q} C(p,q)}$$
(3)

Relaxing the unit norm constraint on $z(\cdot)$ yields a generalized eigenproblem:

$$Wz = \lambda Dz \tag{4}$$

with D and W defined in terms of C and Θ by:

$$D = \operatorname{Diag}(C1_n) \tag{5}$$

$$W = C \bullet e^{i\Theta} \tag{6}$$

where *n* is the number of pixels, 1_n is a column vector of ones, $Diag(\cdot)$ is a matrix with its vector argument on the main diagonal, and \bullet denotes the matrix Hadamard product.

For Θ everywhere zero (W = C), this eigenproblem is identical to the spectral relaxation of Normalized Cuts [29], in which the second and higher eigenvectors encode grouping [29, 2]. With nonzero entries in Θ , the first of the now



Figure 3. Complex affinities for grouping and figure/ground. An angular displacement, corresponding to relative figure/ground or depth ordering, along with a confidence on that displacement, specify pairwise local grouping relationships between pixels. A single complex number encodes confidence as magnitude and displacement as angle from the positive real axis. Four basic interaction types span the space of possible pairwise pixel relationships. **Contiguous region:** Pixels p and q lie in the same region. A vector along the positive real axis represents high confidence on zero relative displacement. Ambiguous boundary: Pixels p and q lie in different regions whose interface admits no cues for discriminating displacement. The shared boundary could be a surface marking or depth discontinuity with either of p or q in front. The origin of the complex plane represents zero confidence on the correct relationship. Figure transition: As boundary convexity tends to indicate foreground, moving from p to q likely transitions from ground to figure. We have high confidence on positive angular displacement. *Ground transition:* In the reverse case, q is ground with respect to p, and the complex representation has negative angle.

complex-valued eigenvectors is nontrivial and its angle encodes rank ordering while the subsequent eigenvectors still encode grouping [22]. We use the same decoding procedure as [22] to read off this information.

Specifically, given eigenvectors, $\{z_0, z_1, ..., z_{m-1}\}$, and corresponding eigenvalues, $\lambda_0 \leq \lambda_1 \leq ... \leq \lambda_{m-1}$, solving Equation 4, $\theta(p) = \arg(z_0(p))$ recovers figure/ground ordering. Treating the eigenvectors as an embedding of pixels into \mathbb{C}^m , distance in this embedding space reveals perceptual grouping. We follow [2, 22] to recover both boundaries and segmentation from the embedding by taking the (spatial) gradient of eigenvectors and applying the watershed transform. This is equivalent to a form of agglomerative clustering in the embedding space, with merging constrained to be between neighbors in the image domain.

A remaining issue, solved by [24], is to avoid circular wrap-around in angular span by guaranteeing that the solution fits within a wedge of the complex plane. It suffices to rescale Θ by $\frac{\pi}{2}(1_n^T|\Theta|1_n)^{-1}$ prior to embedding.

Having chosen Angular Embedding as our inference procedure, it remains to define the pairwise pixel relationships C(p,q) and $\Theta(p,q)$. In the special case of Normalized Cuts, C(p,q) represents a clustering affinity, or confidence on zero separation (in both clustering and figure/ground). For



Figure 4. Generalized affinity. Combining the base cases in Figure 3, we express generalized affinity W as the sum of a binding force acting along the positive real axis, and figure and ground displacement forces acting at angles. In absence of any strong boundary, the binding force dominates, linking pixels together. In presence of a strong and discriminative boundary, either the figure or ground force dominates, triggering displacement. Under conditions of uncertainty, all forces are weak. *Left:* The plot for |W| illustrates total force strength, while the plot for $\angle W$ shows the dominant force. *Right:* Complex-valued W varies smoothly across its configuration space, yet exhibits four distinct modes (binding, figure, ground, uncertain). Smooth transitions occur in the region of uncertainty at the origin.

the more general case, we must also predict non-zero figure/ground separation values and assign them confidences. Let us develop the model in terms of probabilities:

$$e(p) = Pr(p \text{ lies on a boundary})$$
 (7)

$$b(p,q) = Pr(\operatorname{seg}(p) \neq \operatorname{seg}(q))$$
(8)

$$f(p,q) = Pr(\text{figural}(p,q) \mid \text{seg}(p) \neq \text{seg}(q)))$$
(9)

$$g(p,q) = Pr(\text{figural}(q,p) \mid \text{seg}(p) \neq \text{seg}(q)))$$
(10)

where seg(p) is the region (segment) containing pixel p and figural(p,q) means that q is figure with respect to p, according to the true segmentation and figure/ground ordering. b(p,q) is the probability that some boundary separates p and q. f(p,q) and g(p,q) are conditional probabilities of figure and ground, respectively. Note g(p,q) = 1 - f(p,q).

There are three possible transitions between p and q: none (same region), ground \rightarrow figure, and figure \rightarrow ground. Selecting the most likely, the probabilities of erroneously binding p and q into the same region, transitioning to figure, or transitioning to ground are respectively:

$$E_B(p,q) = b(p,q) \tag{11}$$

$$E_F(p,q) = 1 - (1 - e(p))b(p,q)(1 - e(q))f(p,q)$$
(12)

$$E_G(p,q) = 1 - (1 - e(p))b(p,q)(1 - e(q))g(p,q)$$
(13)

where (1 - e(p))b(p,q)(1 - e(q)) is the probability that there is a boundary between p and q, but that neither p nor *q* themselves lie on any boundary. Figure/ground repulsion forces act long-range and across boundaries. We convert to confidence via exponential reweighting:

$$C_B(p,q) = \exp(-E_B(p,q)/\sigma_b) \tag{14}$$

$$C_F(p,q) = \exp(-E_F(p,q)/\sigma_f)$$
(15)

$$C_G(p,q) = \exp(-E_G(p,q)/\sigma_g) \tag{16}$$

where σ_b and $\sigma_f = \sigma_g$ control scaling. Using a fixed angle ϕ for the rotational action of figure/ground transitions $(\Theta(p,q) = \pm \phi)$, we obtain complex-valued affinities:

$$W_B(p,q) = C_B(p,q) \tag{17}$$

$$W_F(p,q) = C_F(p,q) \exp(i\phi) \tag{18}$$

$$W_G(p,q) = C_G(p,q) \exp(-i\phi)$$
(19)

Figure 3 illustrates how W_B (shown in green), W_F (red), and W_G (blue) cover the base cases in the space of pairwise grouping relationships. Combining them into a single energy model (generalized affinity) spans the entire space:

$$W(p,q) = W_B(p,q) + W_F(p,q) + W_G(p,q)$$
(20)

One can regard W(p,q) as a sum of binding, figure transition, and ground transition forces acting between p and q. Figure 4 plots the configuration space of W(p,q) in terms of b(p,q) and f(p,q). As the areas of this space in which each component force is strong do not overlap, W behaves in distinct modes, with a smooth transition between them through the area of weak affinity near the origin.

Learning to predict e(p), b(p,q), and f(p,q) suffices to determine all components of W. For computational efficiency, we predict pairwise relationships between each pixel and its immediate neighbors across multiple spatial scales. This defines a multiscale sparse W. As an adjustment prior to feeding W to the Angular Embedding solver of [23], we enforce Hermitian symmetry by assigning:

$$W \leftarrow (W + W^*)/2 \tag{21}$$

3. Affinity Learning

Supervised training of our system proceeds from a collection of images and associated ground-truth, $\{(I_0, S_0, R_0), (I_1, S_1, R_1), \ldots\}$. Here, I_k is an image defined on domain $\Omega_k \subset \mathbb{N}^2$. $S_k : \Omega_k \to \mathbb{N}$ is a segmentation mapping each pixel to a region id, and $R_k : \Omega_k \to \mathbb{R}$ is an rank ordering of pixels according to figure/ground layering. This data defines ground-truth pairwise relationships:

$$\tilde{b}_k(p,q) = 1 - \delta(S(p) - S(q)) \tag{22}$$

$$\tilde{f}_k(p,q) = (sign(R(q) - R(p)) + 1)/2$$
 (23)

As f(p,q) is a conditional probability (Equation 9), we only generate training examples $\tilde{f}_k(p,q)$ for pairs (p,q) satisfying $\tilde{b}_k(p,q) = 1$.



Figure 5. **Deep Affinity Network.** Our CNN produces an output feature map with spatial resolution matching the input image and whose channels encode pairwise affinity estimates between each pixel and 8 neighbors across 3 scales. Its internal architecture derives from that of previous networks [9, 8, 26] for generating perpixel predictions using multiple receptive fields. At each stage, red labels denote the spatial size ratio of feature map to input image. Blue labels denote the number of channels in each feature map.

In all experiments, we sample pixel pairs (p, q) from a multiscale stencil pattern. For each pixel p, we consider as q each of its 8 immediate neighbors in the pixel grid, across 3 scales (distances of 1, 4, and 16 pixels). The stencil pattern thus consists of 24 neighbors total. We train 48 predictors, $b(\cdot, \cdot)$ and $f(\cdot, \cdot)$ at each of the 24 offsets, for describing the pairwise affinity between a pixel and its neighbors. We derive the predictor $e(\cdot)$ as a local average of $b(\cdot, \cdot)$:

$$e(p) = \frac{1}{8} \sum_{q \in \mathcal{N}_1(p)} b(p,q)$$
 (24)

where $\mathcal{N}_1(p)$ consists of the 8 neighbors to p at fine-scale.

Choosing a CNN to implement these predictors, we regard the problem as mapping an input image to a 48 channel output over the same domain. We adapt prior CNN designs for predicting output quantities at every pixel [9, 8, 26] to our somewhat higher-dimensional prediction task. Specifically, we reuse the basic network design of [26], which first passes a large-scale coarse receptive field through an AlexNet [19]-like subnetwork. It appends this subnetwork's output into a second scale subnetwork acting on a finer receptive field. Figure 5 provides a complete layer diagram. In modifying [26], we increase the size of the penultimate feature map as well as the output dimensionality.

For modularity at training time, we separately train two networks, one for $b(\cdot, \cdot)$ and one for $f(\cdot, \cdot)$, each with the layer architecture of Figure 5. We use modified Caffe [17] for training with log loss between truth \tilde{y} and prediction y applied to each output pixel-wise:

$$\mathcal{L}^{\log}\left(\tilde{y}, y\right) = -\tilde{y}\log(y) - (1 - \tilde{y})\log(1 - y)$$
(25)



Figure 6. Affinity learning for segmentation and figure/ground. Ground-truth assembly (left): Given only ground-truth segmentation [25] and local figure/ground labels in the form of boundary ownership [13], we infer a global ground-truth figure/ground order by running Angular Embedding with pairwise interactions defined by the local ground-truth. Affinity training (right): The ground-truth segmentation serves to train pairwise grouping probability $b(\cdot, \cdot)$, while the globalized ground-truth figure/ground trains $f(\cdot, \cdot)$. Shown are ground-truth training targets \tilde{b} , \tilde{f} , and model predictions b, f, for one component of our stencil: the relationship between pixel p and its neighbor at relative offset d = (-16, 0). Ground-truth \tilde{b} is binary (blue=0, red=1). \tilde{f} is also binary, except pixel pairs in the same region (shown green) are ignored. As f is masked by b at test time, we require only that f(p, q) be correct when b(p, q) is close to 1.

making the total loss for $b(\cdot, \cdot)$:

$$\mathcal{L} = \frac{1}{|\Omega||\mathcal{N}(p)|} \sum_{p \in \Omega} \sum_{q \in \mathcal{N}(p)} \mathcal{L}^{\log}\left(\tilde{b}(p,q), b(p,q)\right) \quad (26)$$

with an analogous loss applied for $f(\cdot, \cdot)$. Here $\mathcal{N}(p)$ denotes all 24 neighbors of p according to the stencil pattern.

Using stochastic gradient descent with random initialization and momentum of 0.9, we train with batch size 32 for 5000 mini-batch iterations. Learning rates for each layer are tuned by hand. We utilize data augmentation in the form of translation and left-right mirroring of examples.

4. Experiments

Training our system for the generic perceptual task of segmentation and figure/ground layering requires a dataset fully annotated in this form. While there appears to be renewed interest in creating large-scale datasets with such annotation [39], none has yet been released. We therefore use the Berkeley segmentation dataset [25] for training. Though it consists of 500 images total, only 200 have been annotated with ground-truth figure/ground [13]. We resplit this subset of 200 images into 150 for training and 50 for testing.

The following subsections detail, how, even with such scarcity of training data, our system achieves substantial improvements in figure/ground quality over prior work.

4.1. BSDS: Ground-truth Figure/Ground

Our model formulation relies on dense labeling of pixel relationships. The BSDS ground-truth provides a dense segmentation in the from of a region map, but only defines local figure/ground relationships between pixels immediately adjacent along a region boundary [13]. We would like to train predictors for long-range figure/ground relationships (our multiscale stencil pattern) in addition to short-range.

Figure 6 illustrates our method for overcoming this limitation. Given perfect (*e.g.* ground-truth) short-range predictions as input, Angular Embedding generates an extremely high-quality global figure/ground estimate. In a real setting, we want robustness by having many estimates of pairwise relations over many scales. Ground-truth short-range connections suffice as they are perfect estimates. We use the globalized ground-truth figure/ground map (column 4 in Figure 6) as our training signal R in Equation 23. The usual ground-truth segmentation serves as S in Equation 22.

4.2. BSDS: Segmentation & Figure/Ground Results

Figure 7 shows results on some examples from our 50 image test set. Compared to the previous attempt [22] to use Angular Embedding as an inference engine for figure/ground, our results are strikingly better. It is visually apparent that our system improves on every single example in terms of figure/ground.

Our segmentation, as measured by boundary quality, is



Figure 7. **Image segmentation and figure/ground results.** We compare our system to ground-truth and the results of Maire [22]. Spectral F/G shows per-pixel figure/ground ordering according to the result of Angular Embedding. The colormap matches Figure 2, with red denoting figure and blue denoting background. Spectral boundaries show soft boundary strength encoded by the eigenvectors. These boundaries generate a hierarchical segmentation [2], one level of which we display in the final column with per-pixel figure/ground averaged over regions. Note the drastic improvement in results over [22]. While [22] reflects a strong lower-region bias for figure, our system learns to use image content and extracts foreground objects. All examples are from our resplit figure/ground test subset of BSDS.

comparable to that of similar systems using spectral clustering for segmentation alone [2]. On the standard boundary precision recall benchmark on BSDS, our spectral boundaries achieve an F-measure of 0.68, identical to that of the spectral component ("spectral Pb") of the gPb boundary detector [2]. Thus, as a segmentation engine our system is on par with the previous best spectral-clustering based systems.

As a system for joint segmentation and figure/ground organization, our system has few competitors. Use of Angular Embedding to solve both problems at once is unique to our system, and [22, 24]. Figure 7 shows an obvious huge jump in figure/ground performance over [22].

4.3. BSDS: Figure/Ground Benchmark

To our knowledge, there is not a well-established benchmarking methodology for dense figure/ground predictions. While [34] propose metrics coupling figure/ground classification accuracy along boundaries to boundary detection performance, we develop a simpler alternative.

Given a per-pixel figure/ground ordering assignment, and a segmentation partitioning an image into regions, we can easily order the regions according to figure/ground layering. Simply assign each region a rank order equal to the mean figure/ground order of its member pixels. For robustness to minor misalignment between the figure/ground assignment and the boundaries of regions in the segmentation, we use median in place of mean.

This transfer procedure serves as a basis for comparing different figure/ground orderings. We transfer them both onto the same segmentation. In particular, given predicted figure/ground ordering $\theta(\cdot)$, ground-truth figure/ground ordering $\tilde{\theta}(\cdot)$, and ground-truth segmentation S, we transfer each of $\theta(\cdot)$ and $\tilde{\theta}(\cdot)$ onto S. This gives two orderings of the regions in S, which we compare according to the following metrics:

- Pairwise region accuracy (R-ACC): For each pair of neighboring regions in *S*, if the ground-truth figure/ground assignment shows them to be in different layers, we test whether the predicted relative ordering of these regions matches the ground-truth relative ordering. That is, we measure accuracy on the classification problem of predicting which region is in front.
- Boundary ownership accuracy (B-ACC): We define the front region as owning the pixels on the common boundary of the region pair and measure the per-pixel accuracy of predicting boundary ownership. This is a reweighting of R-ACC. In R-ACC, all region pairs straddling a ground-truth figure/ground boundary count equally. In B-ACC, their importance is weighted according to length of the boundary.

Segmentation:	Figure/Ground Prediction Accuracy			
Ground-truth	R-ACC	B-ACC	B-ACC-50	B-ACC-25
F/G: Ours	0.62	0.69	0.72	0.73
F/G: Maire [22]	0.56	0.58	0.56	0.56
				<u>.</u>
Segmentation:	Figure/Ground Prediction Accuracy			
Segmentation.	ן гig	ure/Ground	Prediction Ac	curacy
Ours	R-ACC	ure/Ground B-ACC	B-ACC-50	B-ACC-25
Ours F/G: Ours	R-ACC 0.66	B-ACC 0.70	B-ACC-50 0.69	B-ACC-25 0.67
OursF/G: OursF/G: Maire [22]	R-ACC 0.66 0.59	B-ACC 0.70 0.62	B-ACC-50 0.69 0.61	B-ACC-25 0.67 0.58

Table 1. Figure/ground benchmark results. After transferring figure/ground predictions onto either ground-truth (upper table) or our own (lower table) segmentations, we quantify accuracy of local relative relationships. R-ACC is pairwise region accuracy: considering all pairs of neighboring regions, what fraction are correctly ordered by relative figure/ground? B-ACC is boundary ownership accuracy: what fraction of boundary pixels have correct figure ownership assigned? B-ACC-50 and B-ACC-25 restrict measurement to the boundaries of the 50% and 25% most foreground regions (a proxy for foreground objects). Our system dramatically outperforms [22] across all metrics.

Boundary ownership of foreground regions (B-ACC-50, B-ACC-25): Identical to B-ACC, except we only consider boundaries which belong to the foreground-most 50% or 25% of regions in the ground-truth figure/ground ordering of each image. These metrics emphasize the importance of correct predictions for fore-ground objects while ignoring more distant objects.

Note that S need not be the ground-truth segmentation. We can project ground-truth figure/ground onto any segmentation (say, a machine-generated one) and compare to predicted figure/ground projected onto that segmentation.

Table 1 reports a complete comparison of our figure/ground predictions and those of [22] against groundtruth figure/ground on our 50 image test subset of BSDS [25]. We consider both projection onto ground-truth segmentation and onto our own system's segmentation output. For the latter, as our system produces hierarchical segmentation, we use the region partition at a fixed level of the hierarchy, calibrated for optimal boundary F-measure. Figures 8 and 9 provide visual comparison on 10 test images.

Across all metrics, our system significantly outperforms [22]. We achieve 69% and 70% boundary ownership accuracy on ground-truth and automatic segmentation, respectively, compared to 58% and 62% for [22].

4.4. Additional Datasets

Figure 10 demonstrates that our BSDS-trained system captures generally-applicable notions of both segmentation and figure/ground. On both PASCAL VOC [10] and the Weizmann Horse database [5], it generates figure/ground layering that respects scene organization. On the Weizmann examples, though having only been trained for perceptual organization, it behaves like an object detector.



Figure 8. **Figure/ground predictions measured on ground-truth segmentation.** We transfer per-pixel figure/ground predictions (columns 2 through 4 of Figure 7) onto the ground-truth segmentation by taking the median value over each region. For boundaries separating regions with different ground-truth figure/ground layer assignments, we check whether the predicted owner (more figural region) matches the owner according to the ground-truth. The rightmost two columns mark correct boundary ownership predictions in green and errors in red for both the results of Maire's system [22] and our system. Note how we correctly predict ownership of object lower boundaries (rows 6, 8, 10) and improve on small objects (row 7). Table 1 gives quantitative benchmarks.



Figure 9. **Figure/ground predictions measured on our segmentation.** As in Figure 8, we transfer ground-truth figure/ground, Maire's figure/ground predictions [22], and our predictions onto a common segmentation. However, instead of using the ground-truth segmentation, we transfer onto the segmentation generated by our system. The ground-truth figure/ground transferred onto our regions defines the boundary ownership signal against which we judge predictions. Comparing with Figure 8, our boundary ownership predictions are mostly consistent regardless of the segmentation (ground-truth or ours) to which they are applied. However, row 3 shows this is not the case for [22]; here, its predicted correct ownership for the lower boundary relies on averaging out over a large ground-truth background region.





Figure 10. **Cross-domain generalization.** Our system, trained only on the Berkeley segmentation dataset, generalizes to other datasets. Here, we test on images from the PASCAL VOC dataset [10] and the Weizmann Horse database [5]. On PASCAL, our figure/ground reflects object attention and scene layering. On the Weizmann images, our system acts essentially as a horse detection and segmentation algorithm despite having no object-specific training. Its generic understanding of figure/ground suffices to automatically pick out objects.

5. Conclusion

We demonstrate that Angular Embedding, acting on CNN predictions about pairwise pixel relationships, provides a powerful framework for segmentation and figure/ground organization. Our work is the first to formulate a robust interface between these two components. Our results are a dramatic improvement over prior attempts to use spectral methods for figure/ground organization.

References

- [1] P. Arbeláez. Boundary extraction in natural images using ultrametric contour maps. *POCV*, 2006.
- [2] P. Arbeláez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *PAMI*, 2011.
- [3] G. Bertasius, J. Shi, and L. Torresani. DeepEdge: A multiscale bifurcated deep network for top-down contour detection. *CVPR*, 2015.
- [4] G. Bertasius, J. Shi, and L. Torresani. High-for-low and lowfor-high: Efficient boundary detection from deep object features and its applications to high-level vision. *ICCV*, 2015.
- [5] E. Borenstein and S. Ullman. Combined top-down/bottomup segmentation. *PAMI*, 2008.
- [6] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. arXiv:1412.7062, *ICLR*, 2015.
- [7] P. Dollár and C. L. Zitnick. Fast edge detection using structured forests. *PAMI*, 2015.
- [8] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. *CVPR*, 2015.
- [9] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. *NIPS*, 2014.
- [10] M. Everingham, L. van Gool, C. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes (VOC) challenge. *IJCV*, 2010.
- [11] C. Fowlkes, D. Martin, and J. Malik. The Berkeley Segmentation Dataset and Benchmark (BSDB). http://www.cs.berkeley.edu/projects/vision/grouping/segbench/.
- [12] C. Fowlkes, D. Martin, and J. Malik. Learning affinity functions for image segmentation: Combining patch-based and gradient-based approaches. *CVPR*, 2003.
- [13] C. Fowlkes, D. Martin, and J. Malik. Local figure/ground cues are valid for natural images. *Journal of Vision*, 2007.
- [14] Y. Ganin and V. S. Lempitsky. N⁴-fields: Neural network nearest neighbor fields for image transforms. *ACCV*, 2014.
- [15] X. He and A. Yuille. Occlusion boundary detection using pseudo-depth. ECCV, 2010.
- [16] D. Hoiem, A. N. Stein, A. A. Efros, and M. Hebert. Recovering occlusion boundaries from a single image. *ICCV*, 2007.
- [17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv*:1408.5093, 2014.

- [18] Z. Jia, A. Gallagher, Y.-J. Chang, and T. Chen. A learning based framework for depth ordering. *CVPR*, 2012.
- [19] A. Krizhevsky, S.Ilya, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. *NIPS*, 2012.
- [20] G. Lin, C. Shen, I. Reid, and A. van den Hengel. Efficient piecewise training of deep structured models for semantic segmentation. *CVPR*, 2016.
- [21] Y. Lu, W. Zhang, H. Lu, and X. Xue. Salient object detection using concavity context. *ICCV*, 2011.
- [22] M. Maire. Simultaneous segmentation and figure/ground organization using angular embedding. ECCV, 2010.
- [23] M. Maire and S. X. Yu. Progressive multigrid eigensolvers for multiscale spectral segmentation. *ICCV*, 2013.
- [24] M. Maire, S. X. Yu, and P. Perona. Object detection and segmentation from joint embedding of parts and pixels. *ICCV*, 2011.
- [25] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. *ICCV*, 2001.
- [26] T. Narihira, M. Maire, and S. X. Yu. Direct intrinsics: Learning albedo-shading decomposition by convolutional regression. *ICCV*, 2015.
- [27] X. Ren, C. Fowlkes, and J. Malik. Figure/ground assignment in natural images. ECCV, 2006.
- [28] W. Shen, X. Wang, Y. Wang, X. Bai, and Z. Zhang. Deep-Contour: A deep convolutional feature learned by positivesharing loss for contour detection. *CVPR*, 2015.
- [29] J. Shi and J. Malik. Normalized cuts and image segmentation. PAMI, 2000.
- [30] A. N. Stein and M. Hebert. Occlusion boundaries from motion: Low-level detection and mid-level reasoning. *IJCV*, 2009.
- [31] D. Sun, C. Liu, and H. Pfister. Local layering for joint motion estimation and occlusion detection. *CVPR*, 2014.
- [32] P. Sundberg, T. Brox, M. Maire, P. Arbeláez, and J. Malik. Occlusion boundary detection and figure/ground assignment from optical flow. *CVPR*, 2011.
- [33] S. C. Turaga, J. F. Murray, V. Jain, F. Roth, M. Helmstaedter, K. Briggman, W. Denk, and H. S. Seung. Convolutional networks can learn to generate affinity graphs for image segmentation. *Neural Computation*, 2010.
- [34] G. P. Visa and P. Salembier. Precision-recall-classification evaluation framework: Application to depth estimation on single images. *ECCV*, 2014.
- [35] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. Yuille. Joint object and part segmentation using deep learned potentials. *ICCV*, 2015.
- [36] Y. Yang, S. Hallman, D. Ramanan, and C. Fowlkes. Layered object models for image segmentation. *PAMI*, 2011.
- [37] S. X. Yu. Angular embedding: from jarring intensity differences to perceived luminance. CVPR, 2009.
- [38] S. X. Yu. Angular embedding: A robust quadratic criterion. PAMI, 2012.
- [39] Y. Zhu, Y. Tian, D. Mexatas, and P. Dollár. Semantic amodal segmentation. arXiv:1509.01329, 2015.