# DARNet: Deep Active Ray Network for Building Segmentation

Dominic Cheng[1, 2]    Renjie Liao[1, 2, 3]    Sanja Fidler[1, 2, 4]    Raquel Urtasun[1, 2, 3]

University of Toronto[1]    Vector Institute[2]    Uber ATG Toronto[3]    NVIDIA[4]

{dominic, rjliao, fidler}@cs.toronto.edu    urtasun@uber.com

## Abstract

*In this paper, we propose a Deep Active Ray Network (DARNet) for automatic building segmentation. Taking an image as input, it first exploits a deep convolutional neural network (CNN) as the backbone to predict energy maps, which are further utilized to construct an energy function. A polygon-based contour is then evolved via minimizing the energy function, of which the minimum defines the final segmentation. Instead of parameterizing the contour using Euclidean coordinates, we adopt polar coordinates, i.e., rays, which not only prevents self-intersection but also simplifies the design of the energy function. Moreover, we propose a loss function that directly encourages the contours to match building boundaries. Our DARNet is trained end-to-end by back-propagating through the energy minimization and the backbone CNN, which makes the CNN adapt to the dynamics of the contour evolution. Experiments on three building instance segmentation datasets demonstrate our DARNet achieves either state-of-the-art or comparable performances to other competitors.*

## 1. Introduction

The ability to automatically extract building footprints from aerial imagery is an important task in remote sensing. It has many applications such as cartography, urban planning, and humanitarian aid. While maps in well-established urban areas provide precise definitions of building outlines, in more general situations, this information may be neither up-to-date nor available altogether. Such is the case in remote population centers and in dynamic scenarios caused by rapid urban development or natural disasters. These situations motivate the use of automatic building segmentation to form an understanding of an area.

Convolutional neural networks (CNNs) have rapidly established themselves as the de facto standard for tasks of semantic and instance segmentation, as demonstrated by their impressive performance across a variety of datasets [25, 7, 13]. When applied to the task of building segmentation, however, there exists room for improvement. Specifi-
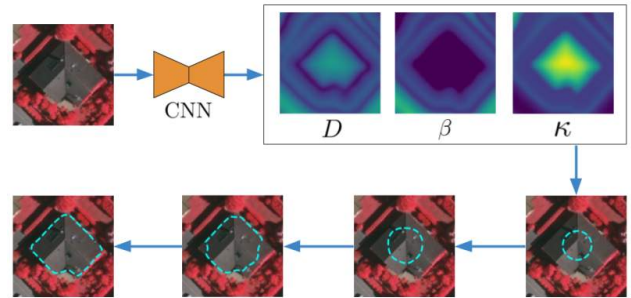


Figure 1: We introduce our DARNet framework, which approaches the problem of instance segmentation by defining a contour using active rays that evolve according to energies parameterized by a CNN. Given an input image, a CNN outputs three maps that define energies. An initialized contour then moves to minimize the energy, yielding an instance segmentation.

cally, as demonstrated in [22], while CNNs are able to generate dense, pixel-based predictions with high recall, they have issues with precise delineation of building boundaries.

Arguably, these boundaries are the most useful features in defining the shape and location of a building. Motivated by this, Marcos *et al*. proposed the deep structured active contours (DSAC) [15] model, which combines the power of deep CNNs with the classic polygon-based active contour model of Kass *et al*. [11]. This fits in particularly well with our prior understanding of buildings, which are generally laid out as relatively simple polygons. DSAC uses a deep CNN to predict an energy landscape on which the active contour, also known as a snake, crawls. When the energy reaches the minimum, the snake stops and the enclosed region is the final predicted segmentation.

Although DSAC's contours exhibit improved coverage compared to CNN based segmentation, they can still be blob-like without strict adherence to building boundaries. Additionally, the representation of contour points with two degrees of freedom presents some challenges. Most notably, it results in extra computational overhead to minimize the proposed energy, and also allows for some contours to exhibit self-intersections. To address these limita-

tions, we introduce the Deep Active Ray Network (DAR-Net), a framework based on a polar representation of active contours, traditionally known as active rays. This ray-based parameterization provides several advantages: 1) contour self-intersections are completely eliminated; 2) it allows us to propose a simpler energy function; 3) the parameterization lends itself to a new loss function that encourages contours to match building boundaries. Our DARNet also exploits a deep CNN as the backbone network for predicting the energy landscape. We train the whole model end-to-end by back-propagating through the energy minimization and the backbone CNN. We compare DARNet against DSAC on three datasets, Vaihingen, Bing Huts, and TorontoCity, and demonstrate its improved effectiveness in producing segmentations that better align with building boundaries.

## 2. Related Work

**Active Contour Models:** First introduced in [11] by the name of *snakes*, active contour models proved to be an extremely popular approach to image segmentation. In these models, an energy is defined as a functional, and its minimization yields a contour that describes a segmentation. The description of this energy is based on intuitive geometric priors and leverages features such as the intensity image and its gradient. Of the myriad works that followed, [5] proposed a balloon force to avoid the tendency for snakes to collapse when initialized far from object edges. [6] reformulated snakes in a polar representation known as *active rays* to reduce the energy optimization from two dimensions to one, and addressed the issue of contour collapse by introducing energies that encouraged circular shapes. Our approach leverages the parameterization of active rays, but we elect to use the curvature term proposed in snakes, as our application of interest does not typically contain circular boundaries. Furthermore, we propose a novel balloon energy that does not involve computing normals at every point in our contour, but rather exploits the properties of a polar parameterization to achieve desired effect in a manner that can be efficiently computed, and fits in seamlessly with our contour inference.

**Deep Active Contour Models:** [18] proposed to combine deep learning with the active contours framework of [20] by having a CNN predict a vector field for a patch around a given contour point to guide it to the closest boundary. However, this method is unable to be learned end-to-end, as CNN learning is separated from the active contour inference. [9, 23] propose to combine a CNN with a level set method [16] in an end-to-end differentiable manner. In contrast to level sets, which define contours implicitly, snakes provide an explicit representation of contour points allowing for the definition of energies based on geometric intuition. [15] uses the snakes framework and replaces the engineered features with ones learned by a CNN.

The problem is posed under the formulation of structured output prediction, and the CNN is trained using a structured support vector machine (SSVM) hinge loss [21] to optimize for intersection-over-union. In contrast, we propose to use an active rays parameterization alongside a largely simplified energy functional. We also propose a loss function that encourages sharper, better aligned contours. Additionally, we back-propagate our loss through the contour evolution, *i.e.*, the energy minimization. It is interesting to note that there are other deep learning based approaches for predicting polygon-based contours. For example, rather than representing a polygon as a discretization of some continuous function, [4, 2] use a recurrent neural network (RNN) to directly predict the polygon vertices in a sequential manner. In [12], the authors predict the polygon or spline outlining the object using a Graph Convolutional Network.

**Building Segmentation:** Current state-of-the-art approaches to building segmentation typically incorporate CNNs in a two stage pipeline: identify instances, and extract polygons. As shown in [22], instances can be extracted from connected components of the semantic mask predicted by a semantic segmentation network [14, 8], or directly predicted with an instance segmentation network [3]. A polygon is then extracted from the instance mask. Because the output space of these approaches are individual pixels, the networks do not reason about the geometry of its predictions, resulting in segmentations that are blob-like around building boundaries. In contrast, the concept of a polygonal output is directly embedded in our model, where we can encourage our outputs to match ground truth shapes.

## 3. Our Approach

In this section, we introduce our DARNet model. Given an input image, our CNN predicts feature maps that define an energy landscape. We place a polygon-based contour, parameterized using polar coordinates, at an initial position on this landscape. This contour then evolves to minimize its energy via gradient descent, and its resting position defines the predicted instance segmentation. We refer the reader to Figure 1 for an illustration of our model. In the subsequent sections, we first describe our parametrization of the contour, highlighting its advantages, such as avoiding self-intersections. We then introduce our energy formulation, in particular our novel balloon energy, and the contour evolution process. Last, we explain how to learn our model in an end-to-end manner.

### 3.1. Contour Representation

Recall that in the active contour (or snake) model a contour point $v$ is represented as

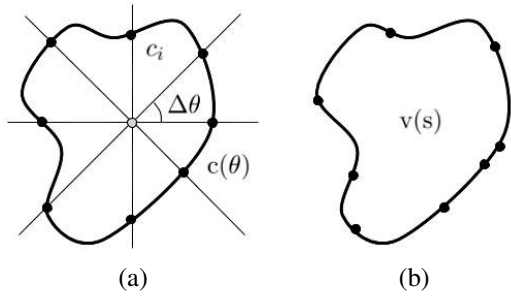$$v(s) = \begin{bmatrix} x_v(s) \\ y_v(s) \end{bmatrix} \tag{1}$$

Figure 2: Contours defined by (a) active rays and (b) snake.



Figure 3: Multiple boundary intersections can occur for (a) non-convex region, (b) self-intersection.

where $s$ denotes the arc length and is generally defined over an interval $s \in [0,1]$. Note that by varying the arc length $s$ from 0 to 1, we obtain the contour. Since this parameterization adopts separate functions for $x$ and $y$ coordinates, the contour point is free to move in any direction, which may cause self-intersection as shown in Figure 3. In contrast, in this paper, we propose to use a parametrization that implicitly avoids self-intersection.

**Active Rays:** Inspired by [6], we parameterize the contour via rays. In particular, we define a contour point $c$ as

$$c(\theta) = \begin{bmatrix} x_c + \rho(\theta)\cos\theta \\ y_c + \rho(\theta)\sin\theta \end{bmatrix} \qquad (2)$$

where $(x_c, y_c)$ define the reference point of the contour, $\rho \geq 0$ is the radius and $\theta$ is the angle tracing out from the x-axis and ranging $[0, 2\pi)$. We assume $(x_c, y_c)$ to be fixed within the interior of the object of interest. To ease the computation, we discretize the contour as a sequence of $L$ points $\{c_i\}_{i=1}^{L}$. The discretization is chosen such that points have equal angular spacing,

$$c_i = \begin{bmatrix} x_c + \rho_i\cos(i\Delta\theta) \\ y_c + \rho_i\sin(i\Delta\theta) \end{bmatrix} \qquad (3)$$

where $\Delta\theta = \frac{2\pi}{L}$ and $\rho_i = \rho(i\Delta\theta)$. The above ray based parameterization is called *active rays*. Importantly, if the region enclosed by the contour forms a convex set, we can guarantee that for any interior reference point, given any angle $\theta$, there is only one corresponding radius $\rho$ based on the following proposition.

**Proposition 1.** *Given a closed convex set $X$, a ray starting from any interior point of $X$ will intersect with the boundary of $X$ once.*

We leave the proof to the supplementary material. If the region is non-convex, a ray may possibly have more than one intersecting point with the boundary. In that case, we pick the one with the minimum distance from the reference point, thus eliminating the possibility that there are
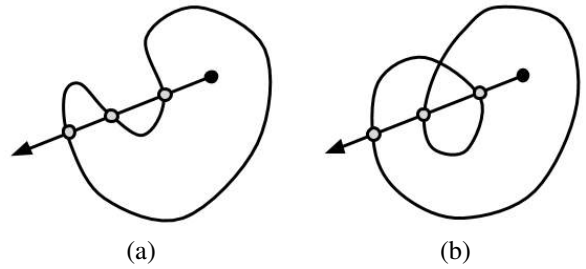
multiple $\rho$ corresponding to the same angle $\theta$. Therefore, compared to snakes, an active rays parameterization avoids self-intersections as shown in Figure 3. Moreover, since we fix the angle at which rays can emanate, the contour points possess an inherent order. Such ordering does not exist for snakes; for example, any cyclic permutation of the snake points produces an identical contour. As we see in Section 3.4, this allows us to naturally use a loss function that encourages our contours to match building boundaries.

**Multiple Sets of Active Rays:** Note that active rays largely preclude contours that enclose non-convex regions. While this is not the dominating case in our application domain, we would like to create a solution that can handle non-convex shapes. Towards this goal, we propose to use multiple sets of active rays, where each set has its own fixed reference point. First, we exploit an instance segmentation to generate a segment over the region of interest (RoI). We use the method in [3] as it tends to under segment the RoIs, thus largely guaranteeing our reference point to lie in the interior. Using this segment, we calculate a distance transform on the mask, and select the location of the largest value as the reference point of our initial contour. If this evolved contour cannot cover the whole segment, we then repeat the process using the distance transform of the uncovered regions, until we cover the whole segment. This is illustrated in Figure 4. The union of evolved contours is used to construct the final prediction.

### 3.2. Contour Energy

We now introduce the formulation of the contour energy functional, of which the minimization yields the final contour. In particular, we encourage the contour to follow boundaries, prefer low-curvature solutions, and expand outwards from a small initialization. With the aforementioned discretization, the overall energy is defined as follows,

$$E(c) = \sum_{i=1}^{L} \left[ E_{\text{data}}(c_i) + E_{\text{curve}}(c_i) + E_{\text{balloon}}(c_i) \right]. \qquad (4)$$
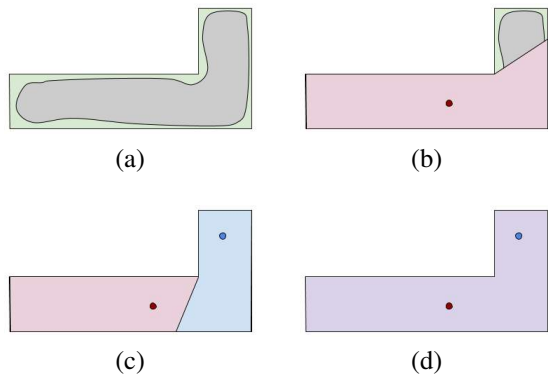
Figure 4: Multiple initialization scheme. (a) Instance segmentation from [3] (gray), and ground truth (green); (b) First initialization and segmentation (red); (c) Initialization from remaining segment (blue); (d) Final output.

Unlike traditional active contour models, we parameterize the energy with the help of a backbone CNN. The hope is that the great power of representation learning of CNN can make the contour evolution more accurate and efficient. We use a CNN architecture based on Dilated Residual Networks [24]. Specifically, we use DRN-D-22, with weights pretrained on ImageNet [19], and append additional learned upsampling layers using transposed convolutions. At the last layer, we predict three outputs that match the input image size, corresponding to three maps, which we denote as $(D, \beta, \kappa)$. We now describe each energy term in detail.

**Data Term:** Given an input image, the data term is defined as

$$E_{\text{data}}(c_i) = D(c_i), \tag{5}$$

where $D$ is a non-negative feature map output by the backbone CNN. Note that $D$ and all subsequently described feature maps are of the same shape as the input image. Since we are minimizing the energy, the contour should seek out places where $D$ is low. Therefore, the CNN should ideally predict low values at the boundaries.

**Curvature Term:** Intuitively, this term models the resistance of the contour towards bending as follows

$$E_{\text{curve}}(c_i) = \beta(c_i)|c_{i+1} - 2c_i + c_{i-1}|^2, \tag{6}$$

where $\beta$ is a non-negative feature map output by the backbone CNN and the squared term is a discrete approximation of the second order derivative of the contour. This term is flexible since the weighting scheme induced by $\beta$ can make the energy locally adaptive. This energy term will force the contour to straighten out wherever the value is high. Our curvature term is simpler compared to the one

in DSAC [15], where in order to prevent snake points from clustering too close together in low-energy areas, they employ an additional membrane term based on the first order derivative of the contour. In contrast, we do not need such a term as our evenly spaced angles guarantee that contour points will not group too closely.

**Balloon Term:** Our balloon term is defined by

$$E_{\text{balloon}}(c_i) = \kappa(c_i)(1 - \rho_i/\rho_{\text{max}}) \tag{7}$$

where $\kappa$ is a non-negative feature map output by the backbone CNN and $\rho_{\text{max}}$ is the maximum radius a contour can reach without crossing the image boundary. The balloon term is designed to propel a given contour point outwards from its reference point, conditioned on the value of the underlying $\kappa$ map. It is necessary due to two reasons. First, the data term may be insufficient to guide the contour towards the boundaries, especially if the contour was initialized far away from them. Second, as noted in [6], the curvature term has an auxiliary effect of shrinking the contour, which can lead to its collapse.

It is interesting to note that DSAC [15] also employs a balloon term which can be expressed using our notation as below,

$$E_{\text{balloon}}^{\text{DSAC}}(v) = - \sum_{(x,y) \in \Omega(v)} \kappa \tag{8}$$

where $\Omega(v)$ denotes the area enclosed by the contour $v$. This term pushes the contour to encapsulate as much of the $\kappa$ map as possible. In our case, due to the active ray parameterization, a contour point can only move along one axis, either towards the reference point or away. Therefore, our balloon term is much simpler as it avoids the need to perform an area integral. Also, as we will see in the following section, our balloon term fits in seamlessly with our inference procedure.

### 3.3. Contour Evolution

Conditioned on the energy terms predicted by the backbone CNN, the second inference step is achieved through energy minimization. To evolve the initial contour towards the optimal one, we first derive the partial derivatives and set them to zero. We then resort to an iterative algorithm to solve the system of partial derivatives.

Specifically, the partial derivatives of the energy terms w.r.t. the contour are derived as below. For the data term,

$$\frac{\partial E_{\text{data}}(c)}{\partial \rho_i} = \frac{\partial D(c_i)}{\partial x} \cos(i\Delta\theta) + \frac{\partial D(c_i)}{\partial y} \sin(i\Delta\theta) \tag{9}$$

where we change the coordinates back to Cartesian to facilitate the computation of derivatives, *e.g.* with a Sobel filter.

For the curvature term, substituting the Cartesian expression of contour points (Equation 3) into the expressions for the energy (Equation 6 and Equation 4), we have,

$$
\begin{aligned}
\frac{\partial E_{\text{curve}}(c)}{\partial \rho_i} &= \frac{\partial}{\partial \rho_i} \sum_{i=0}^{L-1} \beta(c_i)[\rho_{i+1}^2 + 4\rho_i^2 + \rho_{i-1}^2 - \\
&\quad 4\rho_i(\rho_{i+1} + \rho_{i-1})\cos(\Delta\theta) + \\
&\quad 2\rho_{i+1}\rho_{i-1}\cos(2\Delta\theta)] \\
&\approx [2\beta(c_{i+1})\cos(2\Delta\theta)]\rho_{i+2} + \\
&\quad [-4(\beta(c_{i+1}) + \beta(c_{i-1})\cos(\Delta\theta)]\rho_{i+1} + \\
&\quad 2(\beta(c_{i+1}) + 4\beta(c_i) + \beta(c_{i-1}))]\rho_i + \\
&\quad [-4(\beta(c_i) + \beta(c_{i-1}))\cos(\Delta\theta)]\rho_{i-1} + \\
&\quad [2\beta(c_{i-1})\cos(2\Delta\theta)]\rho_{i-2}
\end{aligned}
\tag{10}
$$

where we discard the term arising from the product rule of differentiation as in [15]. We interpret this approximation as treating the $\beta$ map as not varying within the small vicinity of the contour points. Alternatively, we do not wish for the gradient of the $\beta$ map to exert pressure on the contour. Empirically, we found that doing so stabilizes learning, as the network only needs to adjust values of the $\beta$ map without attention to its gradients.

For the balloon term, we use the same approach as the curvature term to obtain the partial derivative,

$$
\frac{\partial E_{\text{balloon}}(c)}{\partial \rho_i} \approx -\frac{\kappa(c_i)}{\rho_{\max}}
\tag{11}
$$

With the above derivation, we have a collection of $L$ partial differential equations w.r.t. individual contour points. We can summarize this system of equations in a compact matrix form,

$$
\frac{\partial E}{\partial \rho} = A\rho + \mathbf{f}
\tag{12}
$$

where $\rho$ is a column vector of size $L$, $A$ is an $L \times L$ cyclic pentadiagonal matrix comprised of $E_{\text{curve}}$ derivatives, and $\mathbf{f}$ is a column vector comprised of $E_{\text{data}}$ and $E_{\text{balloon}}$ derivatives.

This system of partial differential equations can be solved with an iterative method. The approach taken by [11] and [15] is an implicit-explicit method. For the purposes of our implementation, we adopt an explicit method instead, as it avoids the matrix inverse operation. Specifically, the contour evolves according to

$$
\rho^{(t+1)} = \rho^{(t)} - \Delta t \left( A\rho^{(t)} + f \right)
\tag{13}
$$

where $\Delta t$ is a time step hyperparameter. In practice, we found setting $\Delta t$ as $2e^{-4}$ is stable enough for solving the system.
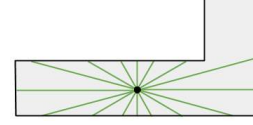


Figure 5: Ground truth rays (green) are defined as the closest distance to the ground truth polygon.

---

**Algorithm 1:** DARNet training algorithm.

---

**Given:** Input image $I$, CNN $F$ with parameters $\omega$, ground truth rays $\rho^{GT}$, initial rays $\rho^{(0)}$, time step hyperparameter $\Delta t$, number of training steps $N$

**for** $j = 1$ **to** $N$ **do**
$\quad (D, \beta, \kappa) = F(I_j; \theta)$
$\quad$ **while** *not converged* **do**
$\quad\quad A, f = $ collection of $\frac{\partial E}{\partial \rho_j^{(t)}}$ using $(D, \beta, \kappa)$
$\quad\quad \rho_j^{(t+1)} = \rho_j^{(t)} - \Delta t(A\rho_j^{(t)} + f)$
$\quad$ **end**
$\quad L = \ell(\rho_j^{(T)}, \rho_j^{GT})$
$\quad$ Compute $\frac{\partial L}{\partial \omega}$ using backpropagation
$\quad$ Update $\omega$ with gradient-based optimization
**end**

---

### 3.4. Learning

Since there exists an explicit ordering of the contour points, we can naturally generate a ground truth active ray and use it to supervise the prediction. Using the same reference point $(x_c, y_c)$ and angle discretization $\Delta\theta$, we cast rays outwards and record the distances at which they intersect the ground truth polygon. In the case of multiple intersections, we take the smallest distance, to prioritize hitting the correct boundaries over increasing coverage. This is illustrated in Figure 5. We use this collection of ground truth distances, $\{\rho_i^{GT}\}_{i=1}^L$, to compute an $L_1$ loss:

$$
\ell(\rho, \rho^{GT}) = \sum_{i=1}^L |\rho_i^{GT} - \rho_i|.
\tag{14}
$$

It differs from the loss employed by DSAC, which used a SSVM hinge loss to optimize for intersection-over-union. Instead, our loss encourages contour points to target building boundaries which is simpler and more efficient.

To allow for gradients to backpropagate to the $D$, $\beta$, and $\kappa$ maps, we interpret the value of a given contour point as a floating point number, and compute the value of a map at that point (*e.g.* $D(c_i)$) using bilinear interpolation from its four adjacent map entries, in a manner similar to what is used in Spatial Transformer Networks [10]. We summarize the learning process in Algorithm 1.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets:** We evaluate DARNet on several building instance segmentation datasets: Vaihingen [1], Bing Huts [15], and TorontoCity [22]. The Vaihingen dataset consists primarily of detached buildings in a town in southern Germany. The original images are $512 \times 512$ at a resolution of 9 cm/pixel. There are 168 buildings in total, which are divided into $100/68$ examples for train/test according to the same split used in [15]. The Bing Huts dataset consists of huts located in a rural area of Tanzania, with an original size of $80 \times 80$ at a resolution of 30 cm/pixel. There are 605 images in total, divided into $335/270$ examples for train/test, again using the same splits. For these two datasets, we further partition the existing training data in an $80\%/20\%$ split to form a validation set, and use this for model selection. We do this for 5 disjoint validation sets, while the test set remains the same. The TorontoCity dataset contains aerial images captured over a dense urban area in Toronto. The images used have a resolution of 10 cm/pixel. The dataset consists of approximately $28,000/12,000$ images for train/test which covers a diverse mixture of buildings, including residential, commercial, and industrial. We divide the training set of TorontoCity into a $80\%/20\%$ split for training and validation respectively.

**Metrics:** We measure performance using intersection-over-union (IoU) averaged over number of instances, weighted coverage, and polygon similarity as in [22]. Additionally, we evaluate the boundary F-score (BoundF) introduced in [17], averaged over thresholds from 1 to 5 pixels, inclusive. For Vaihingen and Bing Huts, we aggregate these metrics over all models selected with the various validation sets, measuring their mean and standard deviation. Lastly, for TorontoCity, we also evaluate the quality of alignment for the predicted boundaries. Specifically, we gather predicted contour pixels that match with the ground truth boundaries, within a threshold of 5 pixels. For these matches, we evaluate the alignment error with respect to the ground truth, which is determined as the cosine similarity between the ground truth boundary and the predicted boundary. We then rank these pixels by their alignment error, and plot the recall at various thresholds of this error.

**Hyper-parameters:** We discretize our contour with $L = 60$ points. For training, we use SGD with momentum, with learning rate $4 \times 10^{-5}$, momentum 0.3, and a batch size of 10. The learning rate decay schedule is explained in supplementary material. We perform 200-step inference as it is found to be sufficient for convergence in practice. We initialize the contour randomly within the ground truth boundary during training. For testing, we initialize in im-
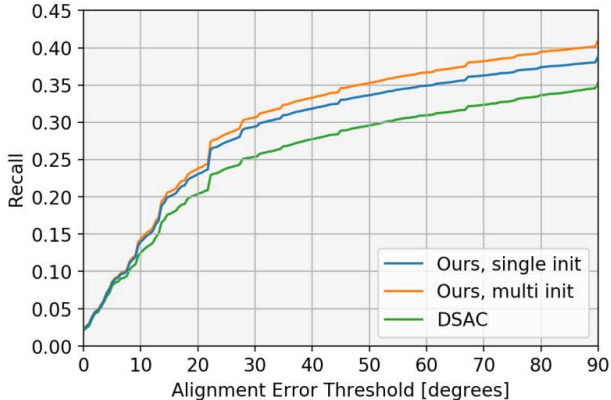


Figure 6: Recall-alignment curve. The recall for all matching boundary pixels (within a 5 pixel threshold of ground truth) is plotted at a given alignment error. Our method exhibits increased recall overall.

age centers for Vaihingen and Bing Huts due to the fact that most buildings of these two datasets are of regular shape. As for TorontoCity, we leverage the deep watershed transform (DWT) [3] to get instance proposals and initialize as described in Section 3.1. Standard data augmentation techniques (random rotation, flipping, scaling, color jitter) are used. We do not adopt common stabilization tricks during inference as they are non-differentiable. Instead, we found using the Euclidean distance transform to pre-train our $D$, $\beta$ and $\kappa$ maps helps stabilize the inference during training. We leave more details of pre-training in the supplementary material.

### 4.2. Results

**Vaihingen and Bing Huts:** We compare against the baseline methods and DSAC [15]. In particular, the baseline exploits a fully convolutional network (FCN) [14] as the backbone to perform semantic segmentation of buildings (interior, boundary, and exterior) and then the post-processing following [15] to obtain the predicted contour. For fair comparison, we also replace the backbone of DSAC and the baseline with ours. We summarize results in Table 1. Compared to the strong FCN baselines, our method exhibits improved performance across the majority of metrics. In particular, the significant improvement on PolySim suggest our segmentations are more geometrically similar. Furthermore, our method significantly outperforms DSAC on all metrics. Even in instances where DSAC exhibits good coverage-based metrics, our method is significantly better at capturing edges. It is interesting to note that substituting our backbone in DSAC does not increase performance, while DSAC's results generally exhibits higher variance, regardless of backbone.

**TorontoCity:** We compare against the semantic segmentation based methods that utilize FCN-8s [14] or

| Method | | Vaihingen | | | | | | | | Bing Huts | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mIoU | | WCov | | PolySim | | BoundF | | mIoU | | WCov | | PolySim | | BoundF | |
| Approach | Backbone | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
| FCN | DSAC's [15] | 81.09 | 0.32 | 81.48 | 0.52 | 68.78 | 0.66 | 64.60 | 0.77 | 69.88 | 0.65 | 73.36 | 0.40 | 64.01 | 0.36 | 30.39 | 0.95 |
| FCN | Ours | 87.27 | 0.50 | 86.89 | 0.40 | 76.41 | 0.77 | **76.84** | 0.48 | 74.54 | 1.22 | **77.55** | 0.96 | 67.98 | 2.20 | 37.77 | 2.69 |
| DSAC [15] | DSAC's [15] | 71.10 | 3.88 | 70.76 | 4.07 | 61.71 | 3.97 | 36.44 | 5.16 | 38.74 | 1.07 | 44.61 | 3.00 | 38.91 | 2.66 | 37.16 | 0.85 |
| DSAC [15] | Ours | 60.37 | 5.97 | 61.12 | 6.22 | 52.96 | 5.89 | 24.34 | 5.73 | 57.23 | 3.89 | 63.09 | 2.25 | 55.43 | 2.20 | 15.98 | 2.62 |
| Ours | Ours | **88.24** | 0.38 | **88.16** | 0.35 | **81.33** | 0.37 | 75.91 | 0.89 | **75.29** | 0.45 | 77.07 | 0.63 | **72.12** | 0.57 | **38.08** | 0.76 |

Table 1: Test set results (mean and standard deviation) on Vaihingen and Bing Huts from models selected across 5 disjoint validation sets. mIoU is averaged over instances. WCov is weighted coverage. PolySim is polygon similarity. BoundF is average boundary-F score at 1 pixel to 5 pixel thresholds.

| Method | mIoU | WCov | PolySim | BoundF |
|---|---|---|---|---|
| FCN-8s [14] | - | 45.6 | **32.3** | - |
| ResNet50 [8] | - | 40.1 | 29.2 | - |
| DWT [3] | - | 52.0 | 24.0 | - |
| DSAC [15] | 60.0 | **58.0** | 27.2 | 25.4 |
| Ours, single init | 57.9 | 52.2 | 23.9 | 29.0 |
| Ours, multi init | **60.1** | 57.5 | 26.8 | **29.6** |

Table 2: Results on TorontoCity.

ResNet50 [8], along with the instance segmentation method that relies on DWT. Additionally, because the commercial and industrial buildings contained in TorontoCity tend to possess complex boundaries, we examine the effects of using one versus multiple initializations described in Section 3.1. We summarize results in Table 2. Our method shows improved performance compared to DSAC on the mIOU. For weighted coverage, which accounts for the area of the building, we achieve comparable performance to DSAC. Our performance on this metric is influenced primarily by larger buildings, as moving from a single initialization to multiple initializations significantly improves the performance. We find that large areas presented by commercial and industrial lots present many ambiguous features such as mechanical penthouses and visible facades due to parallax error. Our single initialization method tends to produce smaller segmentations as it focuses on the boundaries offered by these features. This is alleviated by using multiple initializations. The weighting from larger buildings is also reflected in polygon similarity. The average BoundF metric demonstrates our method is significantly better at capturing building boundaries than DSAC. It is important to note that even our segmentations generated with a single initialization showed improved performance. To delve deeper into the quality of these segmentations, we examine their alignment with respect to the ground truth in Figure 6. We see that, for a given threshold for alignment error, our method exhibits superior recall. Overall, our multiple initialization scheme performs the best, although for lower error thresholds our single initialization is also very competitive.

**Number of Initializations:** In TorontoCity, we found
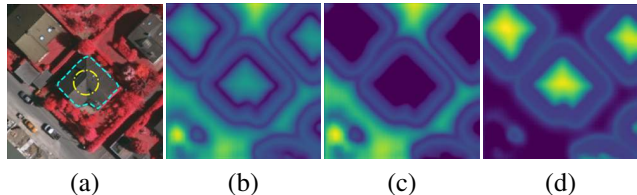


(a)  (b)  (c)  (d)

Figure 7: Energies predicted by our CNN on an example from the Vaihingen test set. (a) Input image, initial contour (yellow), and final contour (cyan). (b) Data term. (c) Curvature term. (d) Balloon term.

that $25.7\%$ of examples required multiple initializations. For these examples, on average 3.71 initializations were required.

**Qualitative Discussion:** We visualize some energies predicted by our CNN in Figure 7. We see the CNN opts to predict a $D$ term that has deep valleys at the building contours. The $\beta$ term adopts small values along the edges to encourage straightness at the boundaries, while the $\kappa$ term acts to propel the contour points from inside. We show additional segmentations in Figure 8. The segmentations produced by our method are generally more adherent to the edges of the buildings. This is especially helpful when buildings are densely packed together, as seen in the TorontoCity results (columns e-f). Additionally, in comparison to DSAC, our parameterization successfully prevents self-intersecting contours (column b, second last row).

**Failure Modes:** The last two rows in Figure 8 demonstrate some weaknesses of our model. In cases where the model is unsure about the extent of one building, it will expand the contour until it meets the edge of another building (column (b), last row). Also, on large commercial lots (column f, last two rows), our method becomes confused by the shapes and features of the buildings.

## 5. Conclusion

In this paper, we presented an approach to building instance segmentation using a combination of active rays with energies predicted by a CNN. The use of a polar representation of contours enables us to predict contours that cannot self-intersect, and to employ a loss function that encourages
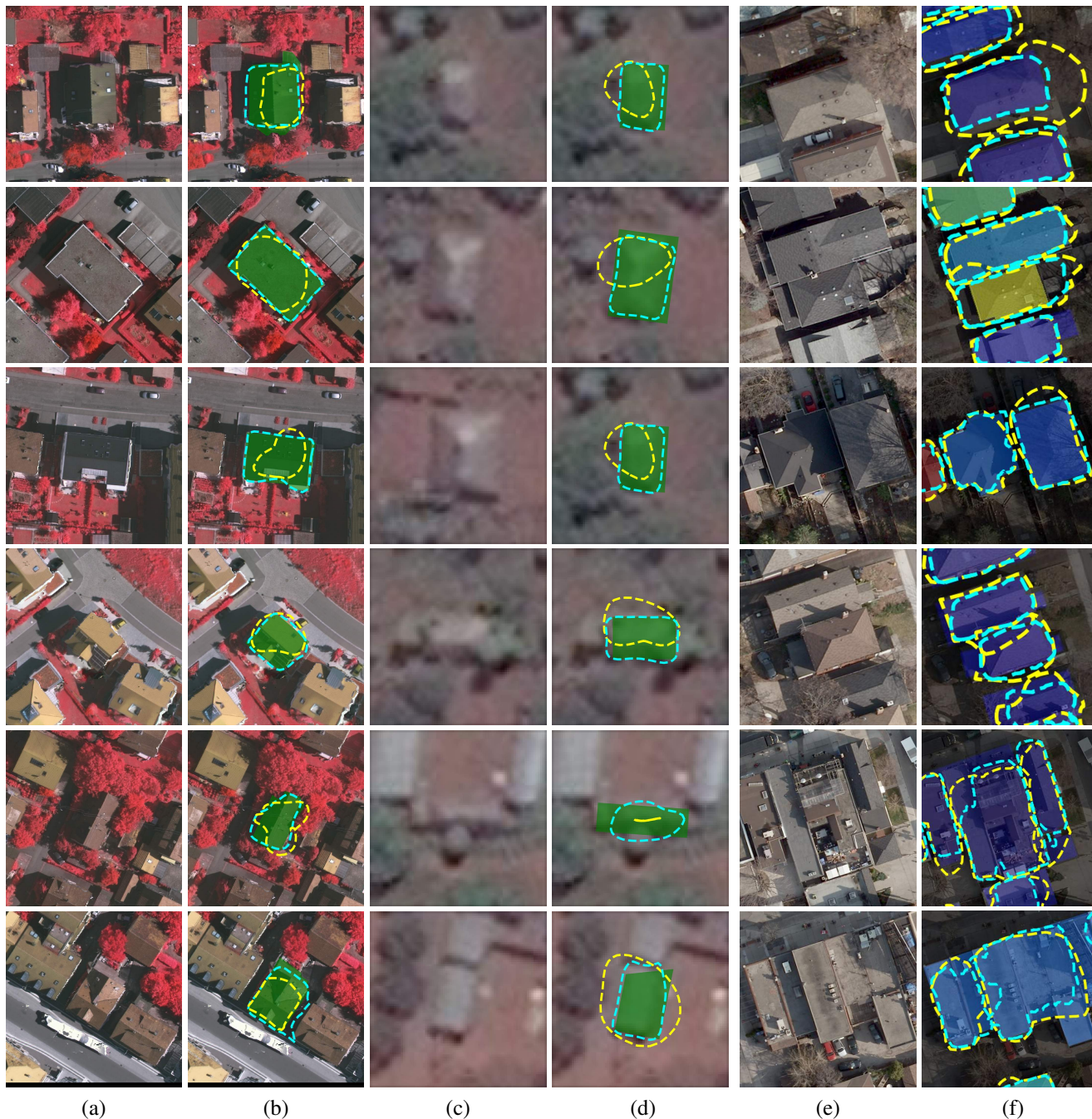
Figure 8: Results on (a-b) Vaihingen, (c-d) Bing Huts, (e-f) TorontoCity. Bottom two rows highlight failure cases. Original image shown in left. On right, our output is shown in cyan; DSAC output in yellow; ground truth is shaded.

our predicted contours to match building boundaries. Furthermore, we demonstrate a method to combine several predictions to generate more complex contours. Comparisons against other state-of-the-art methods on various buliding segmentation datasets demonstrate our method's power in generating segmentations that better capture, and are better aligned with, building boundaries.

## Acknowledgments

# References

[1] International society for photogrammetry and remote sensing, 2d semantic labeling contest. http://www2.isprs.org/commissions/comm3/wg4/semantic-labeling.html. 6

[2] D. Acuna, H. Ling, A. Kar, and S. Fidler. Efficient interactive annotation of segmentation datasets with polygon-rnn++. In *CVPR*, pages 859–868, 2018. 2

[3] M. Bai and R. Urtasun. Deep watershed transform for instance segmentation. In *CVPR*, pages 2858–2866. IEEE, 2017. 2, 3, 4, 6, 7

[4] L. Castrejon, K. Kundu, R. Urtasun, and S. Fidler. Annotating object instances with a polygon-rnn. In *CVPR*, volume 1, page 2, 2017. 2

[5] L. D. Cohen. On active contour models and balloons. *CVGIP: Image understanding*, 53(2):211–218, 1991. 2

[6] J. Denzler, H. Niemann, et al. Active rays: Polar-transformed active contours for real-time contour tracking. *Real-Time Imaging*, 5(3):203–213, 1999. 2, 3, 4

[7] K. He, G. Gkioxari, P. Dollr, and R. Girshick. Mask r-cnn. In *ICCV*, pages 2980–2988, Oct 2017. 1

[8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 2, 7

[9] P. Hu, B. Shuai, J. Liu, and G. Wang. Deep level sets for salient object detection. In *CVPR*, volume 1, page 2, 2017. 2

[10] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *NIPS*, pages 2017–2025, 2015. 5

[11] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active contour models. *IJCV*, 1(4):321–331, 1988. 1, 2, 5

[12] H. Ling, J. Gao, A. Kar, W. Chen, and S. Fidler. Fast interactive object annotation with curve-gcn. In *CVPR*, 2019. 2

[13] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018. 1

[14] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, June 2015. 2, 6, 7

[15] D. Marcos, D. Tuia, B. Kellenberger, L. Zhang, M. Bai, R. Liao, and R. Urtasun. Learning deep structured active contours end-to-end. In *CVPR*, pages 8877–8885, 2018. 1, 2, 4, 5, 6, 7

[16] S. Osher and J. A. Sethian. Fronts propagating with curvature-dependent speed: algorithms based on hamilton-jacobi formulations. *Journal of Computational Physics*, 79(1):12–49, 1988. 2

[17] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Computer Vision and Pattern Recognition*, 2016. 6

[18] C. Rupprecht, E. Huaroc, M. Baust, and N. Navab. Deep active contours. *arXiv preprint arXiv:1607.05074*, 2016. 2

[19] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 4

[20] G. Sundaramoorthi, A. Yezzi, and A. C. Mennucci. Sobolev active contours. *International Journal of Computer Vision*, 73(3):345–366, 2007. 2

[21] I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. Large margin methods for structured and interdependent output variables. *JMLR*, 6:1453–1484, 2005. 2

[22] S. Wang, M. Bai, G. Mattyus, H. Chu, W. Luo, B. Yang, J. Liang, J. Cheverie, S. Fidler, and R. Urtasun. Torontocity: Seeing the world with a million eyes. In *ICCV*, pages 3028–3036, 2017. 1, 2, 6

[23] Z. Wang, D. Acuna, H. Ling, A. Kar, and S. Fidler. Object instance annotation with deep extreme level set evolution. In *CVPR*, 2019. 2

[24] F. Yu, V. Koltun, and T. Funkhouser. Dilated residual networks. In *CVPR*, 2017. 4

[25] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, pages 2881–2890, 2017. 1

# Supplementary Material:
# DARNet: Deep Active Ray Network for Building Segmentation

Dominic Cheng[1, 2]    Renjie Liao[1, 2, 3]    Sanja Fidler[1, 2, 4]    Raquel Urtasun[1, 2, 3]

University of Toronto[1]    Vector Institute[2]    Uber ATG Toronto[3]    NVIDIA[4]

{dominic, rjliao, fidler}@cs.toronto.edu    urtasun@uber.com

## 1. Proof of Proposition

**Proposition 1.** *Given a closed convex set $X$, a ray starting from any interior point of $X$ will intersect with the boundary of $X$ once.*

*Proof.* First, it is straightforward that a ray starting from any interior point of $X$ will intersect with its boundary since otherwise $X$ is not closed. Then we prove the intersection can only happen once by contradiction. We assume a ray starts from an interior point $a$ and intersects with the boundary of $X$ twice at $b$ and $c$. Without loss of generality, we assume $b$ is in between $a$ and $c$. Since $a$ is an interior point, we can find an open ball $A$ which centers at $a$ and $A \in X$. Given any point $\tilde{a} \in A$ other than $a$, we can uniquely determine a line $l$ which crosses $\tilde{a}$ and $c$. Then we can uniquely draw an open ball $B$ which centers at $b$ and has $l$ as the tangent line. Since $b$ is a boundary point of $X$, we can always find a point $\tilde{b} \in B$ such that $\tilde{b} \notin X$. Connecting $c$ and $\tilde{b}$, we can uniquely determine a line $\tilde{l}$. Since $|\angle \tilde{a}cb| \geq |\angle \tilde{b}cb|$, line $\tilde{l}$ will intersect with $A$ for at least once. Denoting any intersection point as $\hat{a}$, we know $\hat{a} \in A \in X$. Since $c \in X$, we know any point in between $c$ and $\hat{a}$, *i.e.*, the convex combination of $c$ and $\hat{a}$, should be in $X$ due to the fact that $X$ is convex. Therefore, $\tilde{b} \in X$ which contradicts. We show the schematic of proof in Fig. 1. $\square$
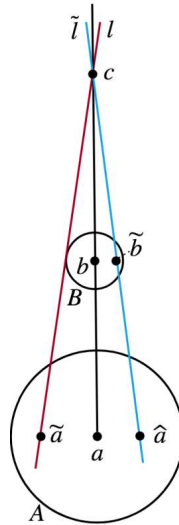


Figure 1: Illustration of proof.

## 2. Contour Inference Details

Our contour inference relies on the following equation,

$$\rho^{(t+1)} = \rho^{(t)} - \Delta t \left( A\rho^{(t)} + f \right) \tag{1}$$

where $\rho^{(t)}$ represents the contour at step $t$, $\Delta t$ is a time step hyper-parameter for solving the system, $A$ and $f$ consist of partial derivatives of the energy w.r.t. $\rho^{(t)}$. Here, we detail the construction of this equation.

**Matrix equation** As mentioned in our paper, the relevant partial derivatives are as follows. For the data term,

$$\frac{\partial E_{\text{data}}(c)}{\partial \rho_i} = \frac{\partial D(c_i)}{\partial x} \cos(i\Delta\theta) + \frac{\partial D(c_i)}{\partial y} \sin(i\Delta\theta) \tag{2}$$

where

$$c_i = \begin{bmatrix} x_c + \rho_i \cos(i\Delta\theta) \\ y_c + \rho_i \sin(i\Delta\theta) \end{bmatrix} \tag{3}$$

For the curvature term,

$$\begin{aligned}
\frac{\partial E_{\text{curve}}(c)}{\partial \rho_i} \approx\ & [2\beta(c_{i+1})\cos(2\Delta\theta)]\rho_{i+2} + \\
& [-4(\beta(c_{i+1}) + \beta(c_{i-1}))\cos(\Delta\theta)]\rho_{i+1} + \\
& 2[\beta(c_{i+1}) + 4\beta(c_i) + \beta(c_{i-1})]\rho_i + \\
& [-4(\beta(c_i) + \beta(c_{i-1}))\cos(\Delta\theta)]\rho_{i-1} + \\
& [2\beta(c_{i-1})\cos(2\Delta\theta)]\rho_{i-2}
\end{aligned} \tag{4}$$

For the balloon term,

$$\frac{\partial E_{\text{balloon}}(c)}{\partial \rho_i} \approx -\frac{\kappa(c_i)}{\rho_{\max}} \tag{5}$$

Combining these into the overall energy, we obtain

$$\begin{aligned}
\frac{\partial E}{\partial \rho_i} \approx\ & [2\beta(c_{i+1})\cos(2\Delta\theta)]\rho_{i+2} + \\
& [-4(\beta(c_{i+1}) + \beta(c_{i-1}))\cos(\Delta\theta)]\rho_{i+1} + \\
& 2[\beta(c_{i+1}) + 4\beta(c_i) + \beta(c_{i-1})]\rho_i + \\
& [-4(\beta(c_i) + \beta(c_{i-1}))\cos(\Delta\theta)]\rho_{i-1} + \\
& [2\beta(c_{i-1})\cos(2\Delta\theta)]\rho_{i-2} + \\
& \frac{\partial D(c_i)}{\partial x}\cos(i\Delta\theta) + \frac{\partial D(c_i)}{\partial y}\sin(i\Delta\theta) - \frac{\kappa(c_i)}{\rho_{\max}}
\end{aligned} \tag{6}$$

We have $L$ such equations, one for each $\rho_i$. In each equation, there are dependencies on the four adjacent entries of $\rho_i$, and $\rho_i$ itself (*i.e.* for entries on the borders, the indices wrap around). This summarizes into matrix form,

$$\frac{\partial E}{\partial \rho} \approx \begin{bmatrix}
c_1 & b_1 & a_1 & 0 & \cdots & 0 & e_1 & d_1 \\
d_2 & c_2 & b_2 & a_2 & 0 & \cdots & 0 & e_2 \\
\vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\
a_{L-1} & 0 & \cdots & 0 & e_{L-1} & d_{L-1} & c_{L-1} & b_{L-1} \\
b_L & a_L & 0 & \cdots & 0 & e_L & d_L & c_L
\end{bmatrix}
\begin{bmatrix} \rho_1 \\ \rho_2 \\ \vdots \\ \rho_{L-1} \\ \rho_L \end{bmatrix}
+ \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{L-1} \\ f_L \end{bmatrix} \tag{7}$$

$$= A\rho + f \tag{8}$$

where

$$a_i = 2\beta(c_{i+1})\cos(2\Delta\theta) \tag{9}$$

$$b_i = -4(\beta(c_{i+1}) + \beta(c_{i-1}))\cos(\Delta\theta) \tag{10}$$

$$c_i = 2(\beta(c_{i+1}) + 4\beta(c_i) + \beta(c_{i-1})) \tag{11}$$

$$d_i = -4(\beta(c_i) + \beta(c_{i-1}))\cos(\Delta\theta) \tag{12}$$

$$e_i = 2\beta(c_{i-1})\cos(2\Delta\theta) \tag{13}$$

$$f_i = \frac{\partial D(c_i)}{\partial x}\cos(i\Delta\theta) + \frac{\partial D(c_i)}{\partial y}\sin(i\Delta\theta) - \frac{\kappa(c_i)}{\rho_{\max}} \tag{14}$$

**Solving the system**  As mentioned in [2], to iteratively minimize the energy, we can solve this system of equations by introducing a time variable such that the solution is found at equilibrium. That is,

$$\frac{\rho^{(t+1)} - \rho^{(t)}}{\Delta t} = -(A\rho + f) \tag{15}$$

where we purposefully leave out the time indices on the right hand side of the equation. [2] then adopts an implicit-explicit method by interpreting $A$ implicitly (so it is associated with $\rho^{(t+1)}$) and $f$ explicitly (so it is associated with $\rho^{(t)}$),

$$\frac{\rho^{(t+1)} - \rho^{(t)}}{\Delta t} = -(A\rho^{(t+1)} + f) \tag{16}$$

$$\rho^{(t+1)} = \left(A + \frac{1}{\Delta t}I\right)^{-1}\left(\frac{1}{\Delta t}\rho^{(t)} - f\right) \tag{17}$$

To avoid the need to perform a batched matrix inverse, we instead adopt an explicit method,

$$\frac{\rho^{(t+1)} - \rho^{(t)}}{\Delta t} = -(A\rho^{(t)} + f) \tag{18}$$

$$\rho^{(t+1)} = \rho^{(t)} - \Delta t(A\rho^{(t)} + f) \tag{19}$$

## 3. CNN Architecture

**Our CNN backbone**  Our CNN backbone uses Dilated Residual Network [5], specifically DRN-D-22 with its last pooling and fully-connected layers removed, as the main method of feature extraction. We append additional upsampling layers to yield output maps that match the input image size. This is illustrated in Figure 2.

**Adaptation to Deep Structured Active Contours (DSAC)**  To implement DSAC [4], we leverage the same architecture in Figure 2, except with 4 outputs corresponding to the energy maps required in that framework. Following the implementation of DSAC, we add a gaussian smoothing layer with kernel size 9 and $\sigma = 2$ to the final output of the data term.

## 4. Hyper-parameters

We train our network using SGD with momentum. We choose a learning rate of $4 \times 10^{-5}$, which halves every $E$ epochs, momentum of 0.3, a weight decay of $1 \times 10^{-5}$, and a batch size of 10. We set $E = 30$ for Vaihingen and Bing Huts, and $E = 1$ for TorontoCity. We train for 100 epochs on Vaihingen and Bing Huts, and 10 epochs on TorontoCity.

To encourage stability in contour inference without using common techniques that are non-differentiable, we pretrain the maps to output values that cause the contour to converge, although not necessarily close to the ground truth rays. Specifically, we leverage the Euclidean distance transform because it possesses some desirable properties. Recall that we wish for $D$ to assign relatively lower values to the building boundaries, such that its gradient near these boundaries can attract contour points towards it. Both of these properties are reflected in the distance transform. For $\beta$, we adopt the distance transform with the building interiors masked out; we wish for contour points to evolve outwards without restriction, and straighten out as it approaches the boundaries. For $\kappa$, we adopt the distance transform with the building exteriors masked out. To pretrain, we take the predictions pr_0 and pr_1 (from Figure 2) and regress to these distance transforms with a smooth $L_1$ loss, using Adam [3] with an initial learning rate of $1 \times 10^{-3}$, which halves every $E$ epochs, and weight decay of $4 \times 10^{-4}$. We set
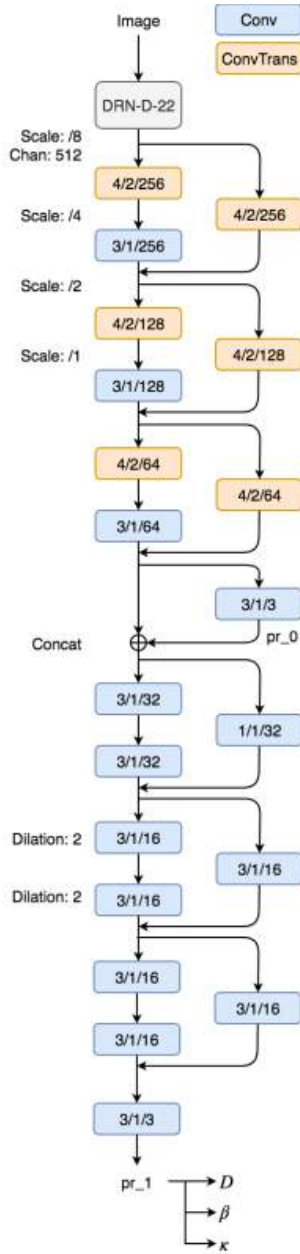
Figure 2: CNN architecture. Convolutions shown in blue; transposed convolutions shown in orange. Layers are annotated as kernel_size/stride/output_channels. Dilation is set to 1 unless stated otherwise. Output sizes are listed as a scale factor from the original input image. All residual blocks are combinations of Convolution-BatchNorm-ReLU in the style of [5], except for pr_0 and pr_1, where BatchNorm [1] is not used.

$E = 50$ for Vaihingen and Bing Huts, and $E = 3$ for TorontoCity. We pretrain for 250 epochs on Vaihingen and Bing Huts, and 10 epochs on TorontoCity. After pretraining, we scale the $\beta$ and $\kappa$ maps by 0.005 and 0.1 respectively so that, during the initial stages of training, the contours do not move too far. We found this procedure increases the stability of training.

## 5. More Visual Examples

We show more examples in Figures 3, 4, 5, and 6.

# References

[1] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015. 4

[2] M. Kass, A. Witkin, and D. Terzopoulos. Snakes: Active contour models. *IJCV*, 1(4):321–331, 1988. 3

[3] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 3

[4] D. Marcos, D. Tuia, B. Kellenberger, L. Zhang, M. Bai, R. Liao, and R. Urtasun. Learning deep structured active contours end-to-end. In *CVPR*, pages 8877–8885, 2018. 3

[5] F. Yu, V. Koltun, and T. Funkhouser. Dilated residual networks. In *CVPR*, 2017. 3, 4

Figure 3: Results on (a-b) Vaihingen, (c-d) Bing Huts, (e-f) TorontoCity. Bottom three rows highlight failure cases. Original image shown in left. On right, our output is shown in cyan; DSAC output in yellow; ground truth is shaded.

Figure 4: Results on (a-b) Vaihingen, (c-d) Bing Huts, (e-f) TorontoCity. Bottom three rows highlight failure cases. Original image shown in left. On right, our output is shown in cyan; DSAC output in yellow; ground truth is shaded.

Figure 5: Results on (a-b) Vaihingen, (c-d) Bing Huts, (e-f) TorontoCity. Bottom three rows highlight failure cases. Original image shown in left. On right, our output is shown in cyan; DSAC output in yellow; ground truth is shaded.

Figure 6: Demonstration of our method (segmentations shown in shaded color) on a large area of the TorontoCity dataset.