

# Annotating Object Instances with a Polygon-RNN

Lluís Castrejón

Kaustav Kundu

Raquel Urtasun

Sanja Fidler

Department of Computer Science

University of Toronto

{castrejon, kkundu, urtasun, fidler}@cs.toronto.edu

## Abstract

We propose an approach for semi-automatic annotation of object instances. While most current methods treat object segmentation as a pixel-labeling problem, we here cast it as a polygon prediction task, mimicking how most current datasets have been annotated. In particular, our approach takes as input an image crop and sequentially produces vertices of the polygon outlining the object. This allows a human annotator to interfere at any time and correct a vertex if needed, producing as accurate segmentation as desired by the annotator. We show that our approach speeds up the annotation process by a factor of 4.7 across all classes in Cityscapes, while achieving 78.4% agreement in IoU with original ground-truth, matching the typical agreement between human annotators. For cars, our speed-up factor is 7.3 for an agreement of 82.2%. We further show generalization capabilities of our approach to unseen datasets.

## 1. Introduction

Semantic image segmentation has been receiving significant attention in the community [5, 17]. With new benchmarks such as Cityscapes [6], object instance segmentation is also gaining steam [14, 24, 34, 21, 29]. Most of the recent approaches are based on neural networks, achieving impressive performance for these tasks [5, 17, 10, 21]. Deep learning approaches are, however, data hungry and their performance is strongly correlated with the amount of available training data. This requires the community to annotate large-scale datasets which is both time consuming and expensive. Our goal in this paper is to make this process faster, while yielding ground-truth as precise as the one available in the current datasets.

There have been several attempts at reducing the dependency on very detailed annotation such as object segmentation masks. In the weakly-supervised setting, approaches aim at learning segmentation models from weak annotation such as image tags or bounding boxes [13, 31, 11]. In [15], the authors rely on scribbles, one on each object, while [1] requires only a single point on the object. While these ap-

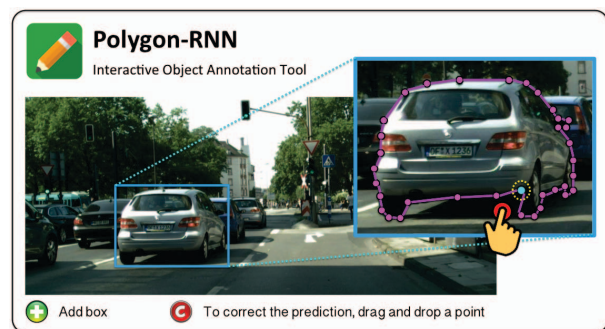


Figure 1. Given a bounding box, we automatically predict the polygon outlining the object instance inside the box, using our Polygon-RNN. Our method is designed to facilitate annotation, and easily incorporates user corrections of points to improve the overall object's polygon. Our method cuts down the number of required annotation clicks by a factor of 4.74.

proaches hold promise, their performance is not yet competitive with fully supervised approaches. Other work exploits easier-to-obtain ground-truth such as bounding boxes, and produces (noisy) labeling inside each box with a GrabCut type of approach [25, 4]. It has been shown that such annotation can serve as useful auxiliary data to train neural segmentation networks [36, 29]. Yet, these segmentations cannot be used as official ground-truth for a benchmark due to its inherent imprecisions.

Most of the large-scale segmentation datasets have been collected by having annotators outline the objects with a polygon [8, 18, 16, 6, 37]. Since typically objects are connected and without holes, polygons provide a way of annotating an object with a relatively small number of clicks, typically around 30 to 40 per object. In this paper, we propose an interactive segmentation method that produces highly accurate and structurally coherent object annotations, and reduces annotation time by a factor of 4.7.

Given a ground-truth bounding box, our method generates a polygon outlining the object instance using a Recurrent Neural Network, which we call Polygon-RNN. Our approach takes as input an image crop and sequentially pro-

duces vertices of the polygon outlining the object. This allows a human annotator to interfere at any time and correct a vertex if needed, producing as accurate segmentations as desired by the annotator. We show that our annotation approach speeds up annotation process by factor of 4.7, while achieving 78.4% agreement with original ground-truth, matching the typical agreement of human annotators. We plan to release our code and create a web-annotation interface running our model at the backend. Please refer to our project page: <http://www.cs.toronto.edu/polyrnn>. We hope this will cut down annotation time and cost of segmentation benchmarks in the future.

## 2. Related Work

Our approach is related to work on semi-automatic image annotation and object instance segmentation.

**Semi-automatic annotation.** There has been significant effort at making pixel-level image labeling faster for the annotators. In [2], the authors used scribbles as seeds to model the appearance of foreground and background, and performed segmentation via graph-cuts by combining appearance cues and a smoothness term [3]. [19] uses multiple scribbles on the object and background and exploits motion cues to annotate an object in a video. Scribbles were also recently used in [15] to train CNNs for semantic image segmentation. GrabCut [25] exploits annotations in the form of 2D bounding boxes, and performs per-pixel labeling with foreground/background models using EM. Building on top of this idea, [23] combined GrabCut with CNN to segment medical images. In [4], the authors exploited 3D bounding boxes and a point cloud to facilitate labeling. A different type of approach has been to exploit multiple bounding boxes and perform co-segmentation [13, 11].

Since most of these approaches define a graphical model at the pixel-level, with the smoothness term as the main relation among pixels, it is hard to incorporate shape priors. These are particularly important in ambiguous regions caused by shadows, image saturation or low-resolution of the object. Furthermore, nothing prevents these models to provide labelings with holes. If the method makes mistakes in outlining the object, the human annotator has a hard and tedious work to correct for such mistakes. Thus, these methods have mainly been used to produce additional, yet noisy training examples, but their output is typically not accurate enough to serve as official ground-truth of a benchmark.

**Annotation tools.** [32] labeled clothing in images by performing annotation at the superpixel-level. This makes the labeling process more efficient, but inherently depends on the superpixel scale and thus typically merges small objects or parts. This issue was somewhat resolved in [22] by labeling videos at multiple superpixel scales.

**Object instance segmentation.** Our work is also related to object instance segmentation. Most of these ap-

proaches [14, 24, 36, 34, 20, 21] operate on the pixel-level, typically exploiting a CNN inside a box or a patch to perform the labeling. Work most related to ours is [35, 28] which aims to produce a polygon around an object. These approaches start by detecting edge fragments and find an optimal cycle that links the edges into a coherent region. In [7], the authors propose a method that produces superpixels in the form of small polygons which they combine into object regions with the aim to label aerial images. In our work, we predict the polygon around the object directly, using a carefully designed RNN.

## 3. Polygon-RNN

Our goal is to create an efficient annotation tool for labeling object instances with polygons. As is typical in an annotation setting, we assume that the user provides the bounding box around the object. Given the image patch inside the box, our method predicts a (closed) polygon outlining the object using a Recurrent Neural Network. We allow the user to correct a predicted vertex of the polygon at any time step if needed, which we integrate in our prediction task.

We parametrize the polygon as a sequence of 2D vertices  $(c_t)_{t \in \mathbb{N}}$ ,  $c \in \mathbb{R}^2$ . We assume the polygon is closed, i.e., there is an edge between any two consecutive vertices, as well as the last and the first vertices. Note that a closed polygon is a cycle and thus has multiple equivalent parametrizations obtained by choosing any of the vertices as the beginning of the sequence, as well as selecting the orientation of the sequence. Here, we fix the polygon to always follow the clockwise orientation, but the starting point of the sequence can be any of the vertices.

Our model is an RNN, that predicts a vertex at every time step. As input in each step of the RNN we use a CNN representation of the image crop, as well as the vertices predicted one and two time steps ago, plus the first point. By explicitly providing information of the past two points we help the RNN to follow a particular orientation of the polygon. On the other hand, the first vertex helps the RNN to decide when to close (finish) the polygon. We train the RNN+CNN model end-to-end. This essentially helps the CNN to be fine-tuned to object boundaries, while the RNN learns to follow these boundaries and exploits its recurrent nature to also encode priors on object shapes. Our model thus returns a structurally coherent representation of the object. We name our model Polygon-RNN.

Figure 2 shows the overview of the model. We next describe each component of the model in more detail.

### 3.1. Model Architecture

We start by providing details on the image representation via a CNN, and then explain the design of the RNN.

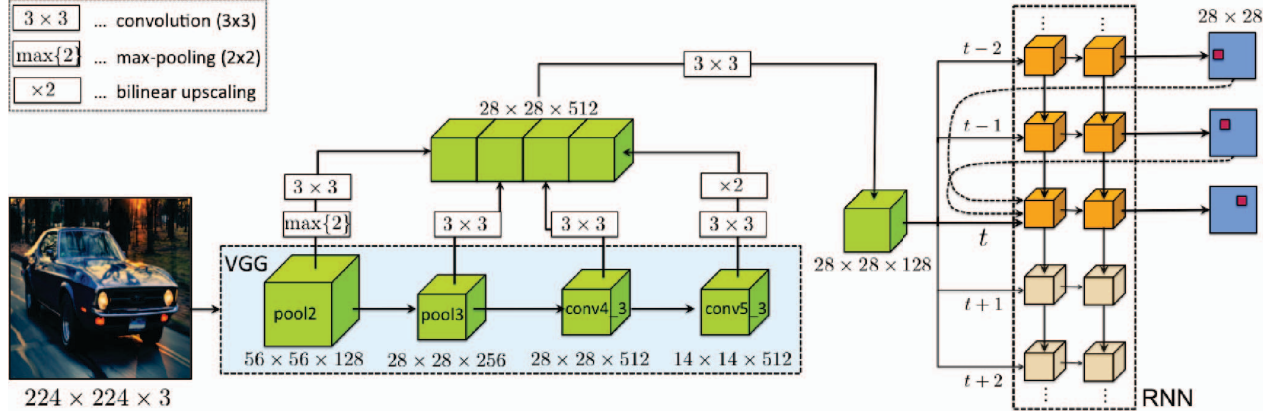


Figure 2. **Our Polygon-RNN model.** At each time step of the RNN-decoder (right), we feed in an image representation using a modified VGG architecture. Our RNN is a two-layer convolutional LSTM with skip-connection from one and two time steps ago. At the output at each time step, we predict the spatial location of the new vertex of the polygon.

### 3.1.1 Image Representation via a CNN with Skip Connections

We adopt the VGG-16 architecture [27] and modify it for the purpose of our task. We first remove the fully connected layers as well as the last max-pooling layer, *pool5*. The output of this modified network has a downsampling factor of 16. We then add additional convolutional layers with skip-connections that fuse information from the previous layers and upscale the output by factor of 2 (downsampling factor of 8 wrt to the original size of the image crop, which is always scaled to  $224 \times 224$ ). This allows the CNN to extract features that contain both low-level information about the edges and corners, as well as semantic information about the object. The latter helps the model to “see” the object, while the former helps it to follow the object’s boundaries.

We employ a similar architecture for the skip-connections as the one in [21]. The design guideline is to first process the features in the skip-layers using another convolutional layer, then concatenate all outputs, and finally process this concatenated feature using another convolutional layer. We employ convolutional filters with a kernel size of  $3 \times 3$ , followed by a ReLU non-linearity. Concatenation layers join the channels of different outputs into a single tensor. Since we use features from multiple skip-layers which have different spatial dimensions, we employ bilinear upsampling or max-pooling in order to get outputs that all have the same spatial resolution. We refer the reader to Fig. 2 for a visualization and further details about the architecture (the CNN is highlighted in green).

### 3.1.2 RNN for Vertex Prediction

An RNN is a powerful representation of time-series data, as it carries more complex information about the history by employing linear and non-linear functions. In our case, we hope the RNN to capture the shape of the object and thus make coherent predictions even in ambiguous cases such as

for example shadows and saturation.

In particular, we employ a Convolutional LSTM [30] in our model, and use it as a decoder. ConvLSTMs operate in 2D, which allows us to preserve the spatial information received from the CNN. Furthermore, a ConvLSTM employs convolutions, thus greatly reducing the number of parameters to be learned compared to using a fully-connected RNN. In its simplest form, a ConvLSTM (single layer) computes the hidden state  $\mathbf{h}_t$  given the input  $\mathbf{x}_t$  according to the following equations:

$$\begin{pmatrix} \mathbf{i}_t \\ \mathbf{f}_t \\ \mathbf{o}_t \\ \mathbf{g}_t \end{pmatrix} = \mathbf{W}_h * \mathbf{h}_{t-1} + \mathbf{W}_x * \mathbf{x}_t + \mathbf{b} \quad (1)$$

$$\mathbf{c}_t = \sigma(\mathbf{f}_t) \odot \mathbf{c}_{t-1} + \sigma(\mathbf{i}_t) \odot \tanh(\mathbf{g}_t)$$

$$\mathbf{h}_t = \sigma(\mathbf{o}_t) \odot \tanh(\mathbf{c}_t)$$

Here  $i, f, o$  denote the input, forget, and output gate,  $h$  is the hidden state and  $c$  is the cell state.  $\sigma$  denotes the sigmoid function,  $\odot$  indicates an element-wise product and  $*$  a convolution.  $W_h$  denotes the hidden-to-state convolution kernel and  $W_x$  the input-to-state convolution kernel.

In particular, we model the polygon with a two-layer ConvLSTM with kernel size of  $3 \times 3$  and 16 channels, which outputs a vertex at each time step. We formulate the vertex prediction as a classification task. Specifically, we represent our output at time step  $t$  as one-hot encoding of a  $D \times D + 1$  grid, where the  $D \times D$  dimensions represent the possible 2D positions of the vertex, and the last dimension corresponds to the end-of-sequence token (i.e., polygon is closed). The position of the vertices are thus quantized to the resolution of the output grid. Let  $y_t$  denote the one-hot encoding of a vertex, output at time step  $t$ .

Our ConvLSTM gets as input a tensor  $\mathbf{x}_t$  at time step  $t$ , that concatenates multiple features: the CNN feature repre-

sensation of the image,  $y_{t-1}$  and  $y_{t-2}$ , i.e., a one-hot encoding of the previous predicted vertex and the vertex predicted from two time steps ago, as well as the one-hot encoding of the first predicted vertex  $y_1$ .

Given two consecutive vertices, the next vertex on the polygon is uniquely defined. However, this is not the case for the first vertex, since any vertex of the polygon can serve as a starting point (polygon is a cycle). We thus treat the starting point as special, and predict it in the following way. We reuse the same architecture of the CNN as in Sec. 3.1.1, but add two layers, each of dimension  $D \times D$ . One branch predicts object boundaries while the other takes as input the output of the boundary-predicting layer as well as the image features and predicts the vertices of the polygon. We treat both, the boundary and vertices as a binary classification problem in each cell in the output grid.

### 3.2. Training

To train our model we use cross-entropy at each time step of the RNN. In order to not over-penalize the incorrect predictions that are close to the ground-truth vertex, we smooth our target distribution at each time step. We assign non-zero probability mass to those locations that are within a distance of 2 in our  $D \times D$  output grid.

We follow the typical training regime where we make predictions at each time step but feed in ground-truth vertex information to the next. We train our model using the Adam optimizer [12] with a batch size  $b = 8$  and an initial learning rate of  $\lambda = 1e - 4$ . We decay the learning rate after 10 epochs by a factor of 10 and use the default values of  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ .

For the task of first vertex prediction, we train another CNN using a multi-task loss. In particular, we use the logistic loss for every location in the grid. As ground-truth for the object boundaries, we draw the edges of the ground-truth polygon, and use the vertices of the polygon as ground-truth for the vertex layer. Our full model takes approximately a day to train on a Nvidia Titan-X GPU.

### 3.3. Inference and Annotators in the Loop

Inference in our model is done by taking the vertex with the highest log-prob at each time step of the RNN. This allows for a simple annotation interface: the annotator can correct the prediction at any time step, and we feed in the corrected vertex to the next time-step of the RNN (instead of the prediction). This puts the model back "on the right track". Typical inference time is 250 ms per object.

### 3.4. Implementation details

We predict the polygon at resolution  $D \times D$ . In our experiments we used  $D = 28$ , corresponding to an 8x down-sampling factor with the input resolution and matching the resolution of the ConvLSTM. We perform polygon simplification with zero error in the quantized grid to eliminate

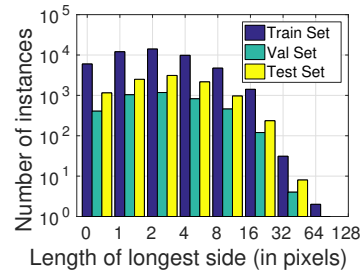


Figure 3. **Distribution of instances across different sizes:** The longest side on the X axis are multiples of 28 pixels

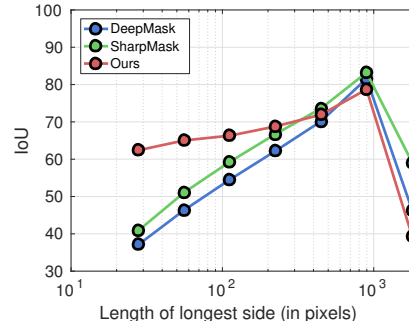


Figure 4. **IoU vs size of instance** comparing different approaches. Here, ours is run in prediction mode.

vertices that lie on a line and to remove multiple vertices that would fall in the same grid position as a result of the quantization process.

We perform three different types of data augmentation: (1) we randomly flip the image crop and the corresponding polygon annotation, (2) we randomly select the amount of context expansion (enlarging the box) between 10% and 20% of the original bounding box and (3) we randomly select the starting vertex of our polygon annotation.

## 4. Results

We evaluate our approach for the task of object instance annotation on the Cityscapes dataset [6], and provide additional results on KITTI [9]. Note that in all our experiments we assume to be given a ground-truth box around the object. Our goal then is to provide a polygon outlining this object as accurately as possible and with minimal number of clicks required from the annotator. We report our performance with the standard IOU measure, as well as the number of vertex corrections of the predicted polygon. A box around the object in principle requires two additional clicks. However, boxes are typically much easier and cheaper to obtain using crowd-sourcing services such as AMT, while for most major segmentation benchmarks, polygons have been collected with high quality (in-house) annotators.

### 4.1. Cityscapes Dataset

We evaluate our approach on the Cityscapes instance segmentation dataset [6]. This dataset has images taken



Split	# Img.	Person	Rider	Car	Truck	Bus	Train	Mbike	Bike
Train	2711	16452	1575	24982	455	352	136	657	3400
Val.	264	1462	180	1962	27	27	32	78	258
Test	500	3395	544	4658	93	98	23	149	1167

Table 1. Number of object instances per class in Cityscapes.

from 27 cities in Germany and neighboring countries. It contains 2975 training, 500 validation and 1525 test images. Since we do not have ground truth instances on the test set, we use an alternative split, where the 500 original validation images form our test set. We then split the original training set and select the images from two cities (Weimar and Zurich) as our validation, while the remaining cities become our training set. The dataset has annotations for eight object categories: *person*, *rider*, *car*, *truck*, *bus*, *train*, *motorcycle* and *bicycle*. The number of instances for each of these classes in our split is shown in Table 1. The Cityscapes dataset has instances with a large variation in their sizes. We show the distribution of instances for different lengths of the longest side of the box, in Fig. 3. We observe a large variance, from 28 pixels to 1792 pixels.

Cityscapes provides instance segmentation ground truth both in terms of a pixel labeling as well as in terms of polygons. In the former, each pixel can correspond to at most one instance, thus representing the visible portion of the object. However, Cityscapes’ polygons typically also capture some occluded parts of an instance, since the annotation tool performed depth ordering of objects to effectively remove the occluded portions [6]. We process the polygons to recreate the layering effect and obtain polygons representing only the visible portions of each object. The average number of vertices from the resulting polygons are shown in Table 2. Since objects can be broken into multiple components due to occlusion, component-wise statistics treats each component as a single example, while instance-wise statistics treats the entire instance as an example. Based on this statistics, we choose a hard limit of 70 time steps for our RNN, taking also GPU memory requirements into account.

**Evaluation Metrics:** We measure two aspects of our predicted annotations. For evaluating their quality, we use the intersection over union (IoU) metric, computed on a per-instance basis, and averaging across all instances. This is a strict measure since the small objects are penalized the same as the large instances. For evaluating the amount of human action required to correct our annotations, we simulate an annotator that corrects a point each time the predicted vertex deviates from the GT vertex more than a threshold. We then report the number of corrections (measured as clicks).

## 4.2. Prediction Mode

We first sanity check the performance of our model without any interaction from the annotator, i.e., we predict the full polygon automatically. We will refer to this setting as the *prediction mode*.

Mode	Car	Truck	Train	Bike	Prsn.	Rider	Mbike	Bus	Avg.
Comp-wise	24.3	27.2	23.6	24.2	27.9	31.6	29.2	26.1	26.8
Inst-wise	31.7	41.7	66.6	40.0	35.0	44.7	45.7	50.8	44.5

Table 2. Average number of vertices in polygon annotations for different object classes in Cityscapes.

**Baselines:** We use the recently proposed DeepMask [20] and SharpMask [21] as state-of-the-art baselines. Given an input image patch, DeepMask uses a CNN to output a pixel labeling of an object, and does so agnostic to the class. Sharpmask extends Deepmask by clever upsampling of the output to obtain the labeling at a much higher resolution (160 vs 56). Note that in their original approach, [20, 21] exhaustively sample patches at different scales over the entire image. Here, we use ground-truth boxes when reporting performance for their approach. Further, DeepMask and SharpMask use a 50 layer ResNet [10] architecture, which has been trained on the COCO [16] dataset. We fine-tune this network on our Cityscapes split in two steps. In the first step, we fine-tune the feed-forward ResNet architecture for 150 epochs, followed by fine-tuning the weights for the Sharpmask’s upsampling layers, for 70 epochs. This two step process is in the same spirit as that suggested in the paper. Note that while these two approaches perform well in labeling the pixels, their output cannot easily be corrected by an annotator in cases when mistakes occur. This is in contrast to our approach, which efficiently integrates a human in the loop in order to get high quality annotations.

We use two additional baselines, SquareBox and Dilation10. SquareBox is a simple baseline where the full box is labeled as the object. Instead of taking the tight-fit box, we reduce the dimensions of the box, keeping the same aspect ratio. Based on the validation set, we get the best results by choosing 80% of the original box. If an instance has multiple components, we fit a box for each individual component as opposed to using the full box. This baseline mimics the scenario, in which the object is modeled simply as a box rather than a polygon. For the Dilation10 baseline, we use the segmentation results from [33], which was trained on the Cityscapes segmentation dataset. For each bounding box, we consider the pixels belonging to the respective object category as the instance mask.

**Quantitative Results:** We report the IoU metric in Table 3. We outperform the baselines in 6 out of 8 categories, as well as in the average across all classes. We perform particularly well in *car*, *person*, and *rider*, outperforming Sharpmask by 12%, 7%, and 6%, respectively. This is particularly impressive since Sharpmask uses a more powerful ResNet architecture (we use VGG).

**Effect of object size:** In Fig. 4, we see how our model performs w.r.t baselines on different instance sizes. For small instances our model performs significantly better than the baselines. For larger objects, the baselines have an ad-

Model	Bicycle	Bus	Person	Train	Truck	Motorcycle	Car	Rider	Mean
Square Box	35.41	53.44	26.36	39.34	54.75	39.47	46.04	26.09	40.11
Dilation10	46.80	48.35	49.37	44.18	35.71	26.97	61.49	38.21	43.89
DeepMask [20]	47.19	69.82	47.93	62.20	63.15	47.47	61.64	52.20	56.45
SharpMask [20]	52.08	<b>73.02</b>	53.63	<b>64.06</b>	65.49	51.92	65.17	56.32	60.21
Ours	<b>52.13</b>	69.53	<b>63.94</b>	53.74	<b>68.03</b>	<b>52.07</b>	<b>71.17</b>	<b>60.58</b>	<b>61.40</b>

Table 3. **Performance** (IoU in %) on all the Cityscapes classes **without the annotator in the loop**.

Threshold	Num. Clicks	Mean IOU
1	15.79	84.74
2	11.77	81.43
3	9.39	78.40
4	7.86	75.79

Table 4. **Annotator in the loop**: Average number of corrections per instance and IoU, computed across all classes. Threshold indicates chessboard distance to the closest GT vertex.

vantage due to larger output resolution. This effect is most notable for classes such as bus and train, in which our model obtains lower IOU compared to the baselines.

### 4.3. Annotator in the loop

The main advantage of our model is that it allows a human annotator to easily interfere if a mistake occurs. In particular, at each RNN time step, the annotator has the possibility to correct a misplaced vertex. The correction is fed to the model at the next time step replacing the model’s prediction, effectively helping the model to get back to the right track. Our goal is to obtain high quality annotations while minimizing annotation time.

We analyze how many clicks are needed to obtain different levels of segmentation accuracy. We perform such analysis by simulating an annotator: we correct a prediction if it deviates from the ground truth vertex by a certain distance. Distances are computed at the model output resolution using the *chessboard* metric. In our experiments we compare the corrected predictions using distance thresholds  $T \in [1, 2, 3, 4]$ . In Table 4 we show the resulting IoU given different thresholds on the distance. We can observe a trade-off between the number of corrections and these metrics.

To put our results in perspective, we hired an experienced, high-quality annotator. We asked the annotator to annotate all car (including van) instances in 10 randomly selected Cityscapes images from our validation split. We perform two experiments: in the first experiment, the annotator is asked to annotate objects by free-viewing of the full image. In the second experiment, we crop the image patches using the Cityscapes boxes, and place a blue dot on the instance to disambiguate annotation. We take a crop with 15% of context around the box and scale it to size 224x224. The annotator used the LabelMe tool [26] for annotation.

In Table 5 we report the IoU achieved by the human annotator as well as the mean number of clicks per instance

Method	Num. Clicks	IoU	Annot. Speed-Up
Cityscapes GT	33.56	100	-
Ann. full image	79.94	69.5	-
Ann. crops	96.09	78.6	-
Ours (Automatic)	0	73.3	No ann.
Ours (T=1)	9.3	87.7	x3.61
Ours (T=2)	6.6	85.7	x5.11
Ours (T=3)	5.6	84.0	x6.01
Ours (T=4)	4.6	82.2	x7.31

Table 5. **Our model vs Annotator Agreement**: We hired a highly trained annotator to label *car* instances on additional 10 images (101 instances). We report IoU agreement with Cityscapes GT, and report polygon statistics. We compare our approach with the agreement between the human annotators.

in each experiment. We can observe that the agreement achieved in IoU is 69.5% in the free-viewing regime, and 78.60% when shown the crops (our regime). This number sheds light on what we are to expect from automatic methods in general, and points to some ambiguities in the task. It also indicates that benchmarks should collect multiple annotations of images to reduce the variations and biases across the annotators. We hope our approach will make such data collection feasible and affordable.

Notice that our model achieves a higher agreement (82%) by requiring only 4.6 clicks on average, which is a factor of 7.3 speed-up in annotation time. Even at agreement as high as 87.7, the annotation speed-up factor is still 3.6. This showcases the effectiveness of our model as an annotation tool. For all the categories in Cityscapes and following the same procedure, we require only 9.39 clicks on average to obtain 78.40% IoU agreement, obtaining a speed-up factor of 4.74.

**Comparison with Grabcut.** We also compare the performance of our approach with another semi automatic method on a set of 54 randomly chosen instances. We used the OpenCV implementation of Grabcut [25] for this experiment. On average, using Grabcut the annotators needed 42.2s and 17.5 clicks per instance, and obtained an average of 70.7% IoU agreement with the Cityscapes GT. On the same set of images, our model achieves IoUs ranging from 79.7% to 85.8%, with 5.0 clicks (T=4) to 9.6 clicks (T=1), respectively. Our expert human annotator needed 87.6 clicks to obtain an IoU of 77.6% (without using any

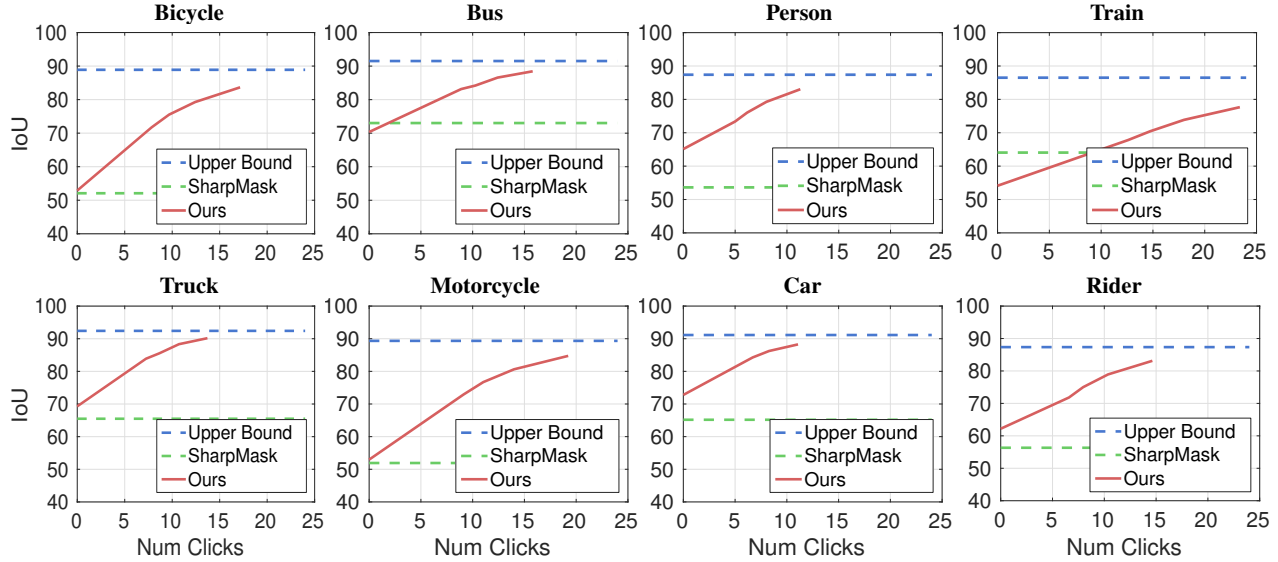


Figure 5. **Annotator in the loop:** We show IoU as a function of the number of clicks/corrections.

Method	# of Clicks	IOU
DeepMask [20]	-	78.3
SharpMask [21]	-	78.8
Beat the MTurkers [4]	0	73.9
Ours (Automatic)	0	74.22
Ours (T=1)	11.83	89.43
Ours (T=2)	8.54	87.51
Ours (T=3)	6.83	85.70
Ours (T=4)	5.84	84.11

Table 6. Car annotation results on the **KITTI** dataset.

semi automatic tool). Since our model requires much less human intervention than [25] (5 vs 17.5 clicks) and requires comparable inference time per click, we expect that in a real world scenario our method would be much faster.

**Qualitative Results:** In Fig. 6 we show examples of images annotated with our method. We remind the reader, that this labeling is obtained by exploiting the GT bounding boxes. In particular, we here show the predictions obtained without any corrections (0 clicks). Our model is able to correctly segment instances with a variety of shapes and sizes. For large instances the quantization error introduced by the output resolution of our model becomes apparent. Increasing the output resolution is subject of ongoing work. The main challenges are memory considerations as well as challenges due to longer sequences (polygons have more vertices) that the network would need to predict.

In Fig. 7 we compare annotations of example instances more closely by zooming in on each object. We can inspect the agreement between the GT annotation and our in-house annotator, as well as the quality of the predictions obtained by PolygonRNN with and without corrections.

#### 4.4. Annotating KITTI Instances

We also evaluate how well our model that was trained on Cityscapes generalizes to an unseen dataset. We use KITTI for this experiment, which has 741 annotated instances provided by [4]. We report the results in Table 6. The object instances in KITTI are usually larger than those found in Cityscapes, making Deepmask and SharpMask perform very similarly. Note that [4], also a method for semi-automatic annotation, exploited Velodyne point clouds to perform their labeling, which puts it with an unfair advantage. Our model is further penalized by its lower resolution output. Still, their performance is lower than our fully automatic approach. With only 5.84 clicks on mean per instance our models achieves an IOU comparable to the human annotation agreement, thus reducing the annotation cost.

#### 5. Conclusion

In this paper we proposed an approach to facilitate annotation of object instances. Our Polygon-RNN predicts a polygon outlining an object, and easily incorporates corrections from an annotator in the loop. We show annotation speed-up of factor 4.74 while achieving the same annotation agreement as that between human annotators. The main advantage of our approach is that it produces structurally plausible annotations of objects, and allows us to achieve a desired annotation accuracy by requiring only a few clicks by the annotator. Additional experiments show that our approach generalizes across different datasets, thus showcasing its power as a generic annotation tool.

#### Acknowledgement

We acknowledge the support from NSERC, and thank Relu Patrascu for infrastructure support. L.C. was supported by a La Caixa Fellowship.





Figure 6. **Qualitative results in prediction mode:** We show polygons for all classes in the original image. Note that our approach uses GT boxes as input. **(left)** we show the GT labeling of the image, **(right)** we show our polygons without any human intervention. The GT images contain 38, 12, 28 and 16 instances, and required 985, 308, 580 and 338 clicks respectively from their Cityscapes annotators.

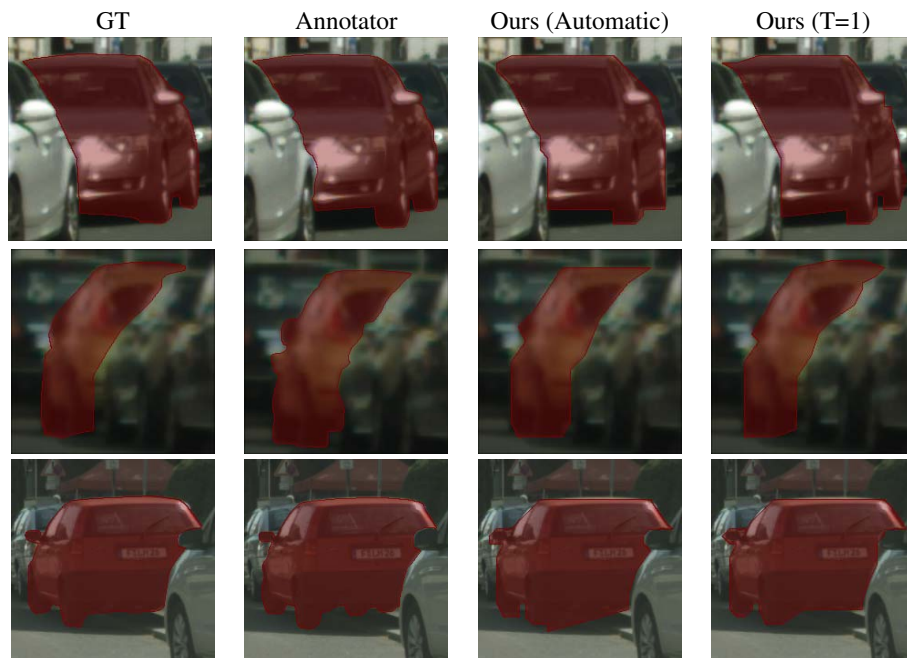


Figure 7. We look at a few instances in more detail. In the **first column** we show the GT annotation, while in the **second column** we show the polygons from the in-house annotator. We observe that these segmentations are high quality but differ in uncertain areas such as the base of the car. In the **third column** we show the PolygonRNN prediction without human intervention. Finally, in the **fourth column** we show a corrected prediction. We can observe that the segmentation is refined to better outline the car mirrors or wheels.



## References

- [1] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei. What's the point: Semantic segmentation with point supervision. *arXiv:1506.02106*, 2016. [1](#)
- [2] Y. Boykov and M.-P. Jolly. Interactive graph cuts for optimal boundary & region segmentation of objects in nd images. In *ICCV*, 2001. [2](#)
- [3] Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *PAMI*, 26(9):1124–1137, 2004. [2](#)
- [4] L.-C. Chen, S. Fidler, A. Yuille, and R. Urtasun. Beat the mTurkers: Automatic image labeling from weak 3d supervision. In *CVPR*, 2014. [1](#), [2](#), [7](#)
- [5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. [1](#)
- [6] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. [1](#), [4](#), [5](#)
- [7] L. Duan and F. Lafarge. Towards large-scale city reconstruction from satellites. In *ECCV*, 2016. [2](#)
- [8] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. [1](#)
- [9] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *CVPR*, 2012. [4](#)
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. [1](#), [5](#)
- [11] S. D. Jain and K. Grauman. Active image segmentation propagation. In *CVPR*, 2016. [1](#), [2](#)
- [12] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [4](#)
- [13] D. Kuettel, M. Guillaumin, and V. Ferrari. Segmentation propagation in imagenet. In *ECCV*, 2012. [1](#), [2](#)
- [14] K. Li, B. Hariharan, and J. Malik. Iterative instance segmentation. In *CVPR*, 2016. [1](#), [2](#)
- [15] D. Lin, J. Dai, J. Jia, K. He, and J. Sun. Scribble-sup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, 2016. [1](#), [2](#)
- [16] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. [1](#), [5](#)
- [17] J. Long, E. Shelhamer, and T. Darrell. Fully Convolutional Networks for Semantic Segmentation. *arXiv:1411.4038*, 2014. [1](#)
- [18] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. *CVPR*, 2014. [1](#)
- [19] N. S. Nagaraja, F. R. Schmidt, and T. Brox. Video segmentation with just a few strokes. In *ICCV*, 2015. [2](#)
- [20] P. O. Pinheiro, R. Collobert, and P. Dollár. Learning to segment object candidates. In *NIPS*, pages 1990–1998, 2015. [2](#), [5](#), [6](#), [7](#)
- [21] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learning to refine object segments. *ECCV 2016*, 2016. [1](#), [2](#), [3](#), [5](#), [7](#)
- [22] J. Pont-Tuset, M. A. F. Guiu, and A. Smolic. Semi-automatic video object segmentation by advanced manipulation of segmentation hierarchies. In *Intl Workshop on Content-Based Multimedia Indexing*, 2015. [2](#)
- [23] M. Rajchl, M. C. Lee, O. Oktay, K. Kamnitsas, J. Passerat-Palmbach, W. Bai, M. Damodaram, M. A. Rutherford, J. V. Hajnal, B. Kainz, and D. Rueckert. Deepcut: Object segmentation from bounding box annotations using convolutional neural networks. In *arXiv:1605.07866*, 2016. [2](#)
- [24] B. Romera-Paredes and P. H. S. Torr. Recurrent instance segmentation. In *arXiv:1511.08250*, 2015. [1](#), [2](#)
- [25] C. Rother, V. Kolmogorov, and A. Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *SIGGRAPH*, 2004. [1](#), [2](#), [6](#), [7](#)
- [26] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman. Labelme: a database and web-based tool for image annotation. *International journal of computer vision*, 77(1-3):157–173, 2008. [6](#)
- [27] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [3](#)
- [28] X. Sun, C. M. Christoudias, and P. Fua. Free-shape polygonal object localization. In *ECCV*, 2014. [2](#)
- [29] J. Uhrig, M. Cordts, U. Franke, and T. Brox. Pixel-level encoding and depth layering for instance-level semantic labeling. In *arXiv:1604.05096*, 2016. [1](#)
- [30] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-k. Wong, and W.-c. Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *NIPS*, pages 802–810, 2015. [3](#)
- [31] J. Xu, A. Schwing, and R. Urtasun. Tell me what you see and i will show you where it is. In *CVPR*, 2014. [1](#)
- [32] K. Yamaguchi, M. H. Kiapour, L. E. Ortiz, and T. L. Berg. Parsing clothing in fashion photographs. In *CVPR*, 2012. [2](#)
- [33] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. [5](#)
- [34] Z. Zhang, S. Fidler, and R. Urtasun. Instance-level segmentation for autonomous driving with deep densely connected mrfs. In *CVPR*, 2016. [1](#), [2](#)
- [35] Z. Zhang, S. Fidler, J. W. Waggoner, Y. Cao, J. M. Siskind, S. Dickinson, and S. Wang. Super-edge grouping for object localization by combining appearance and shape information. In *CVPR*, 2012. [2](#)
- [36] Z. Zhang, A. Schwing, S. Fidler, and R. Urtasun. Monocular object instance segmentation and depth ordering with cnns. In *ICCV*, 2015. [1](#), [2](#)
- [37] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Semantic understanding of scenes through ade20k dataset. In *arXiv:1608.05442*, 2016. [1](#)