

# ShapeMask: Learning to Segment Novel Objects by Refining Shape Priors

Weicheng Kuo<sup>1</sup>, Anelia Angelova<sup>1</sup>, Jitendra Malik<sup>2</sup>, Tsung-Yi Lin<sup>1</sup>  
<sup>1</sup> Google Brain    <sup>2</sup> University of California, Berkeley

<sup>1</sup>{weicheng, anelia, tsungyi}@google.com, <sup>2</sup> malik@eecs.berkeley.edu

## Abstract

Instance segmentation aims to detect and segment individual objects in a scene. Most existing methods rely on precise mask annotations of every category. However, it is difficult and costly to segment objects in novel categories because a large number of mask annotations is required. We introduce ShapeMask, which learns the intermediate concept of object shape to address the problem of generalization in instance segmentation to novel categories. ShapeMask starts with a bounding box detection and gradually refines it by first estimating the shape of the detected object through a collection of shape priors. Next, ShapeMask refines the coarse shape into an instance level mask by learning instance embeddings. The shape priors provide a strong cue for object-like prediction, and the instance embeddings model the instance specific appearance information. ShapeMask significantly outperforms the state-of-the-art by 6.4 and 3.8 AP when learning across categories, and obtains competitive performance in the fully supervised setting. It is also robust to inaccurate detections, decreased model capacity, and small training data. Moreover, it runs efficiently with 150ms inference time on a GPU and trains within 11 hours on TPUs. With a larger backbone model, ShapeMask increases the gap with state-of-the-art to 9.4 and 6.2 AP across categories. Code will be publicly available at: <https://sites.google.com/view/shapemask/home>.

## 1. Introduction

Instance segmentation is the task of providing pixel-level classification of objects and identifying individual objects as separate entities. It is fundamental to applications such as autonomous driving or robot manipulation [8, 44], since segmenting individual objects could help autonomous agents’ planning and decision making. The community has made great headway on this task recently [38, 39, 19, 17, 36, 10, 26, 2, 33, 21, 23]. However, these approaches require precise *pixelwise* supervision for every category. The need for annotation limits instance segmentation to a small slice of visual world that we have dense

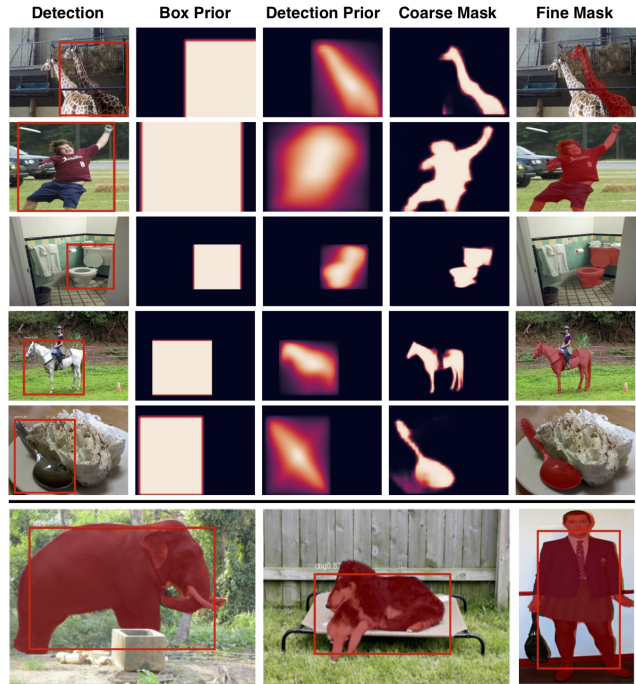


Figure 1: ShapeMask instance segmentation is designed to learn the shape of objects by refining object shape priors. Starting from a bounding box (leftmost column), the shape is progressively refined in our algorithm until reaching the final mask (rightmost column). The bounding box is only needed to approximately localize the object of interest and is not required to be accurate (bottom row).

annotations for. But how can instance segmentation generalize better to novel categories?

Existing instance segmentation algorithms can be categorized into two major approaches: detection-based [38, 39, 19, 17, 10] and grouping-based [35, 1, 36, 26, 2, 33]. To generalize to novel categories, detection-based approaches can use *class-agnostic* training which treats all categories as one foreground category. For example, previous works perform figure-ground segmentation inside a box region without distinguishing object classes [38, 39]. Although class agnostic learning can be readily applied to novel categories, there still exists a clear gap compared to the fully super-

vised setup [21, 38]. On the other hand, the grouping-based approaches learn instance specific cues such as pixel affinity for grouping each instance. Although the grouping stage is inherently class-agnostic and suitable for novel categories, most algorithms still rely on semantic segmentation [1, 35, 2] to provide class information, which requires pixelwise annotation of every class. Whether detection or grouping-based, generalization to novel categories remains an open challenge.

We propose to improve generalization in instance segmentation (Figure 1) by introducing intermediate representations [29, 43, 11], and instance-specific grouping-based learning [40, 23]. Consider Figure 2. Most detection-based approaches use boxes as the intermediate representation for objects (see middle column) which do not contain information of object pose and shape. On the contrary, shapes are more informative (see right column) and have been used by numerous algorithms to help object segmentation [1, 47, 20, 7, 46]. As the pixels of novel objects may appear very different, we hypothesize that shapes can be leveraged to improve generalization as well. Intuitively speaking, learning shapes helps because objects of different categories often share similar shapes, e.g., horse and zebra, orange and apple, fork and spoon. On the other hand, grouping-based learning causes the model to learn “which pixels belong to the same object” and may generalize well by learning appropriate appearance embeddings. For example, even if the model has never seen an orange before, it can still segment it by grouping pixels with similar appearance.

Motivated by these observations, we propose a new instance segmentation algorithm “ShapeMask” to address the generalization problem. Figure 1 illustrates how ShapeMask starts with a box detection, and gradually refines it into a fine mask by learning intermediate shapes. Given a detection, ShapeMask first represents it as a uniform box prior. Then ShapeMask finds the shape priors which best indicate the location, scale and rough shape of the object to fit the box (detection prior). Finally, ShapeMask decodes the coarse mask by a fully convolutional network and refines it by its own instance embedding. The idea behind refinement is similar to grouping approaches. To generalize to novel categories, we simply use class agnostic training for ShapeMask without the need of transfer learning. A natural by-product of learning shapes as soft priors is that ShapeMask can produce masks outside the detection box similar to [18] and unlike [19, 10] which apply feature cropping.

Experiments on COCO show that ShapeMask significantly outperforms the state-of-the-art transfer learning approach [21] in the cross-category setup. In fact, ShapeMask can outperform the state-of-the-art using only 1% of the labeled data. We also qualitatively show that ShapeMask is able to segment many novel object classes in a robotics environment different from the COCO dataset. In the fully su-

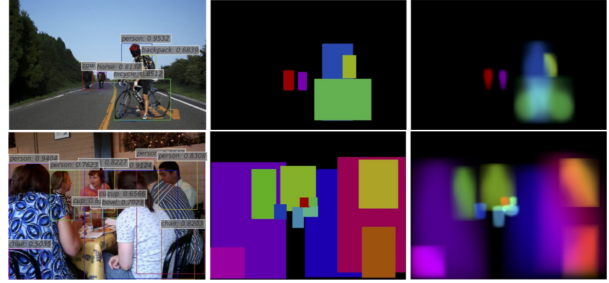


Figure 2: Illustration of objects in uniform box priors vs. shape priors. Every row contains: (left) input image plus detections, (center) box priors, (right) shape priors. Shape priors represent objects with much richer details than boxes.

pervised instance segmentation setting, ShapeMask is competitive with state-of-the-art techniques while training multiple times faster and testing at 150-200ms per image, because it runs seamlessly across hardware accelerators such as TPUs[22, 16] and GPUs to maximize performance.

## 2. Related Work

Instance segmentation can be categorized into two major approaches: *detection-based* and *grouping-based* approaches. The detection-based approaches [17, 10, 18, 28, 6, 19, 38, 39] first detect the bounding box for each object instance and predict the segmentation mask inside the region cropped by the detected box. This approach has been the dominant approach to achieve state-of-the-art performance in instance segmentation datasets like COCO[32] and Cityscapes [8]. The grouping-based approaches [26, 2, 4, 35, 33, 1, 25] view the instance segmentation as a bottom-up grouping problem. They do not assign region of interest for each object instance. Instead, they produce pixelwise predictions of cues such as directional vectors [33], pairwise affinity [35], watershed energy [2], and semantic classes, and then group object instances from the cues in the post-processing stage. In addition to grouping, some object segmentation works have simultaneously used *shape priors* as unaries in probabilistic framework [1, 47, 20], augmented proposals [7], or as top-down prior to help grouping [46, 24, 3]. Classical instance segmentation approaches are mostly grouping-based and work well on unseen data [42, 40]. For example, MCG [40] generates quality masks by normalized cut on the contour pyramid computed from low level cues. So far, grouping-based approaches have not been shown to outperform detection-based methods on the challenging COCO dataset.

Recently, [37, 48, 23, 21] study instance segmentation algorithms that can generalize to categories without mask annotations. [23] leverages the idea that given a bounding box for target object, we can obtain pseudo mask label from a grouping-based segmentation algorithm like Grab-Cut [42]. [37] studies open-set instance segmentation by

using a boundary detector followed by grouping, while [48] learns instance segmentation from image-level supervision by deep activation. Although effective, these approaches do not take advantage of *existing* instance mask labels to achieve better performance.

In this paper, we focus on the *partially supervised* instance segmentation problem [21], as opposed to the weakly-supervised setting [23, 48]. The main idea is to build a large scale instance segmentation model by leveraging large datasets with bounding box annotations e.g. [27], and smaller ones with detailed mask annotations e.g. [32]. More specifically, the setup is that only box labels (not mask labels) are available for a subset of categories at training time. The model is required to perform instance segmentation on these categories at test time. Mask<sup>X</sup> R-CNN [21] tackles the problem by learning to predict weights of mask segmentation branch from the box detection branch. This transfer learning approach shows significant improvement over class-agnostic training, but there still exists a clear gap with the fully supervised system.

### 3. Method

In the following sections, we discuss the set of modules that successively refine object box detections into accurate instance masks.

#### 3.1. Shape Recognition

**Shape priors:** We obtain a set of shape bases from a collection of mask annotations in order to succinctly represent the canonical poses and shapes of each class. These bases are called “shape priors”. The intuition is that when the approximate shape is selected early on in the algorithm, the subsequent instance segmentation becomes much more informed than a box (see also Figure 2). In order to obtain shape priors, we run k-means to find  $K$  centroids of all instance masks for each class in the training set. We resize all mask annotations to a canonical size  $32 \times 32$  before clustering. In the class specific setting, the total number of shape priors is  $C \times K$ , where  $C$  is the number of classes (e.g.  $K = 20$ ). In the class agnostic setting, we group all classes as one and have  $K$  shape priors in total (e.g.,  $K = 100$ ). We define the set of shape priors as  $H = \{S_1, S_2, \dots, S_K\}$ . Figure 3 visualizes example selected shape priors per class for the COCO dataset. We can see the objects have diverse within- and between-class appearance. In class-agnostic setting, clustering yields similarly diverse shape priors.

**Shape estimation:** Starting with a box detection, we first represent it as a binary heatmap  $B$ , i.e.  $b \in \{0, 1\}, \forall b \in B$ . The purpose of this stage is to estimate a more informative detection prior  $S_{prior}$  from  $B$  (see Figure 4). To achieve this, we estimate the target object shape by selecting similar shape priors from the knowledge base  $H$ . Unlike existing methods [6, 19] which view shape prediction as a per-pixel

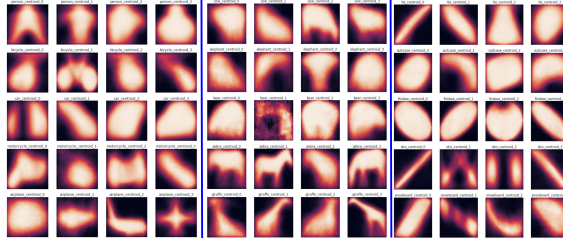


Figure 3: Shape priors obtained by clustering mask labels in the training set. Each prior is a cluster centroid of an object category.

classification problem, we learn to combine similar shapes from  $H$  to form predictions.

Figure 4 illustrates the entire process. First, we pool features inside the bounding box  $B$  on the feature map  $X$ , to obtain an embedding  $x_{box}$  representing the object instance:

$$x_{box} = \frac{1}{|B|} \sum_{(i,j) \in B} X_{(i,j)} \quad (1)$$

The instance shape embedding  $x_{box}$  is then used to recognize similar shapes in the knowledge base  $H$ . The shape priors are the bases used to reconstruct the target object shape inside the bounding box. The predicted object shape  $S$  is a weighted sum of shape priors  $\{S_1, S_2, \dots, S_K\}$ , where the weights are predicted by applying a linear layer  $\phi$  to  $x_{box}$  followed by a softmax function to normalize weights over  $K$ ,  $w_k = \text{softmax}(\phi_k(x_{box}))$

$$S = \sum_{k=1}^K w_k S_k \quad (2)$$

The predicted shape  $S$  is then resized and fitted into the detection box  $B$  to create a smooth heatmap, which we call “detection prior”  $S_{prior}$  (as shown in Figure 4). During training, we apply pixel-wise mean square error (MSE) loss on the detection prior  $S_{prior}$  against the ground-truth mask  $S_{gt}$  to learn the parameters in  $\phi$ .

The approach simplifies the instance segmentation problem by first solving the shape recognition problem. It incorporates the strong prior that the primitive object shapes only have a few modes. This regularizes the output space of the model and prevents it from predicting implausible shapes, e.g., “broken pieces”. By adding such structure to the model, we observe improved generalization to novel categories. We speculate this is because many novel objects share similar shapes with the labeled ones.

#### 3.2. Coarse Mask Prediction

Given the detection prior  $S_{prior}$  from the previous section, the goal of this stage is to obtain a coarse instance mask  $S_{coarse}$  (Figure 5). First, we use a function  $g$  to embed  $S_{prior}$  into the same feature dimension as the image

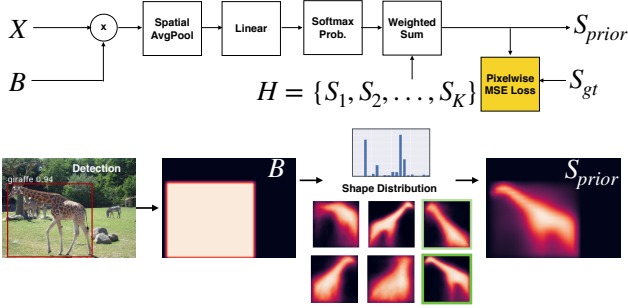


Figure 4: **Shape Estimation.** Given a box detection, we refine the box into an initial estimate of shape  $S_{prior}$  by linearly combining prior shapes  $S_1, S_2, \dots, S_k$ . Our model learns to predict the shape prior distribution to minimize reconstruction error.

features  $X$ , where  $g$  is a  $1 \times 1$  convolution layer. Then we sum them into a prior conditioned feature map  $X_{prior}$ :

$$X_{prior} = X + g(S_{prior}) \quad (3)$$

$X_{prior}$  now contains information from both image features and the detection prior which guides the network to predict object-like segmentation mask. A coarse instance mask  $S_{coarse}$  is decoded by applying a function  $f$  to  $X_{prior}$ , which consists of four convolution layers in our case. This is similar to the mask decoder design in [19], but the difference is we use detection prior  $S_{prior}$  to guide decoding. Pixel-wise cross entropy loss is applied to the predicted mask  $S_{coarse}$  to learn the parameters in the mask decoder:

### 3.3. Shape Refinement by Instance Embedding

Although the coarse segmentation mask  $S_{coarse}$  provides strong cues for possible object shapes, it does not leverage the instance-specific information encoded by the image features. As opposed to previous stages that aim to extract rough shape estimates, the goal of this stage is to refine  $S_{coarse}$  into a detailed final mask  $S_{fine}$  (Figure 6).

Similar to the instance shape embedding  $x_{box}$  in Sec. 3.1, we can pool the instance mask embedding by the refined shape prior to obtain more accurate instance representations  $x_{mask}$ . Given a predicted coarse mask  $S_{coarse}$ , we compute the instance embedding  $x_{mask}$  of the target object by pooling features inside the coarse mask:

$$x_{mask} = \frac{1}{|S_{coarse}|} \sum_{(i,j) \in S_{coarse}} X_{prior(i,j)} \quad (4)$$

We then center the image features  $X_{prior}$  from Equation 3 by subtracting the instance embedding  $x_{mask}$  at all pixel locations:

$$X_{inst(i,j)} = X_{prior(i,j)} - x_{mask} \quad (5)$$

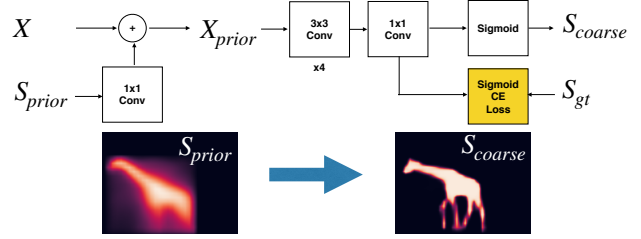


Figure 5: **Coarse Mask Prediction.** We fuse  $S_{prior}$  with the image features  $X$  to obtain prior conditioned features  $X_{prior}$ , from which we decode a coarse shape  $S_{coarse}$ .

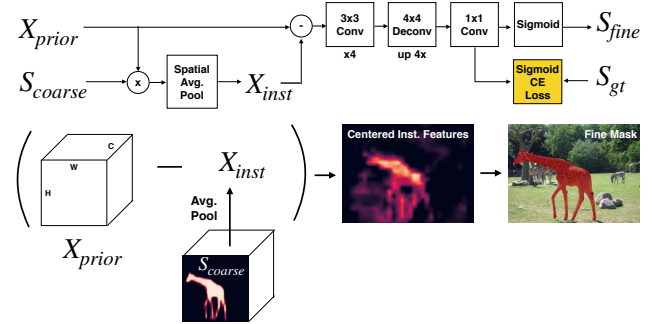


Figure 6: **Shape Refinement.** Starting from  $X_{prior}$  and  $S_{coarse}$ , we first compute the instance embedding  $X_{inst}$  by average pooling the features within  $S_{coarse}$ . Then we subtract  $X_{inst}$  from  $X_{prior}$  before decoding the final mask. We show the low-dimensional PCA projection of the “Centered Instance features” for the purpose of visualization.

This operation can be viewed as conditioning the image features by the target instance. The idea is to encourage the model to learn simple, low-dimensional features to represent object instances. To obtain the fine mask  $S_{fine}$ , we add the mask decoding branch which has the same architecture as described in Section 3.2 with one additional up-sampling layer to enhance the output resolution. Same as before, pixelwise cross entropy loss is used to learn the fine mask  $S_{fine}$  from the groundtruth mask  $S_{gt}$ .

Note that the  $S_{gt}$  here is of higher resolution than before due to the upsampling of  $S_{fine}$ .

### 3.4. Generalization to Class Agnostic Learning

To generalize to novel categories, we adopt class-agnostic learning in ShapeMask. We follow the setup in [21], the box branch outputs box detections with confidence scores for all classes and the mask branch predicts a foreground mask given a box without knowing the class. For generating shape priors  $S_1, S_2, \dots, S_k$ , we combine instance masks from *all* classes together and run k-means with a larger  $K$  than the class-specific setting. This allows us more capacity to capture the diverse modes of shapes across all categories. At inference time, we treat any novel object as part of this one foreground category during shape

estimation and mask prediction stages. The capability to generalize well across categories makes ShapeMask also a *class-agnostic* algorithm, although its performance in the class-specific setting remains competitive among the best techniques.

### 3.5. Implementation Details

We adopt RetinaNet<sup>1</sup> [31] to generate bounding box detections for ShapeMask. Unlike [19, 6] which sample masks from the object proposals, we directly sample 8 groundtruth masks and their associated boxes per image to jitter them for training. Given a bounding box, we assign the box to a feature level in feature pyramid [30] by its longest side and take a fixed-size feature patch centered on the box. More details on the detector, training and feature cropping processes can be found in the Supp. Materials.

## 4. Experiments

**Experimental setup:** We report the performance of ShapeMask on the COCO dataset [32]. We adopt well established protocol in the literature for evaluation [14, 41, 19, 10, 28, 9, 6] by reporting standard COCO metrics AP, AP50, AP75, and AP for small/medium/large objects. Unless specified otherwise, mask AP is reported instead of box AP. We additionally compare the training and inference times, so as to demonstrate the performance/complexity tradeoff.

### 4.1. Generalization to Novel Categories

We first demonstrate the state-of-the-art ability of ShapeMask to generalize across classes and datasets. Such generalization capability shows ShapeMask can work well on a larger part of the visual world than other approaches which require strong pixelwise labeling for every category.

Partially Supervised Instance Segmentation is the task of performing instance segmentation on a subset of categories for which no masks are provided during training. The model is trained on these categories with only box annotations, and on other categories with both box and mask annotations. The experiments are set up following the previous work [21]. We split the COCO categories into “voc” vs. “non-voc”. The voc categories are those also present in PASCAL VOC [12]. At training time, our models have access to the bounding boxes of all categories, but the masks only come from either voc or non-voc categories. The performance upper bounds are set by the oracle models that have access to masks from all categories. In this section, our training set is COCO train2017 and the comparison with other methods is done on val2017 non-voc/voc categories following previous work [21].

**Main results:** We achieve substantially better results than the state-of-the-art methods as shown in Table 1. All benchmark experiments use ResNet-101 network with feature pyramid connections [30]. Using the same FPN backbone, ShapeMask outperforms the state-of-the-art method Mask<sup>X</sup> R-CNN [21] by 6.4 AP on voc to non-voc transfer, and 3.8 AP on non-voc to voc transfer. The gap relative to the oracle upper-bound is 4.8 and 7.6 AP for ShapeMask, compared to the 10.6 and 9.6 AP of Mask<sup>X</sup> R-CNN (lower is better). By adding a stronger feature pyramid from [13], we outperform Mask<sup>X</sup> R-CNN by 9.4 and 6.2 AP. This shows that ShapeMask can take advantage of large backbone model. We also observe that ShapeMask clearly outperforms the baseline class agnostic Mask R-CNN reported in [21] or our own Mask R-CNN implementation. These results provide strong evidence that ShapeMask can better generalize to categories without mask annotations.

Figure 7 visualizes the outputs of ShapeMask in the partially supervised setting. ShapeMask is able to segment many objects well despite not having seen any example mask of the same category during training. The mask branch was trained on voc, tested on non-voc categories and vice versa. By using shape prior and instance embedding, ShapeMask is able to predict complete object-looking shapes in cases where the pixelwise prediction approaches like Mask R-CNN tend to predict broken pieces.

**Generalization with less data:** To study the generalization capabilities of ShapeMask with less training data, we train class agnostic ShapeMask and Mask R-CNN on voc and test on non-voc categories using only 1/1, 1/2, until 1/1000 of the data. To mimic the realistic setting of having less labeled data, we subsample the training set by their image id. Figure 8 shows that ShapeMask generalizes well to unseen categories even down to 1/1000 of the training data. In fact, using just 1/100 of the training data, ShapeMask still outperforms the state-of-the-art Mask<sup>X</sup> R-CNN trained on the whole data by 2.0 AP.

**Generalization to robotics data:** We further demonstrate the ShapeMask algorithm in an out-of-sample scenario, by testing it on object instance segmentation for robotics grasping (Figure 9). This dataset contains many objects not defined in the COCO vocabulary, therefore serving as a good testbed to assess the generalization of ShapeMask. The dataset comes with bounding box annotations on office objects and architectural structures, but *without any instance mask annotation*. The model is *only trained on COCO* and not on this data. To isolate the task of instance segmentation from detection, we feed in groundtruth boxes and evaluate only on segmentation task. As seen, ShapeMask generalizes well to many categories not present in the training data. This shows our approach is particularly useful in settings where the agent will encounter objects beyond the pixelwise annotated vocabulary.

<sup>1</sup><https://github.com/tensorflow/tpu/tree/master/models/official/retinanet>

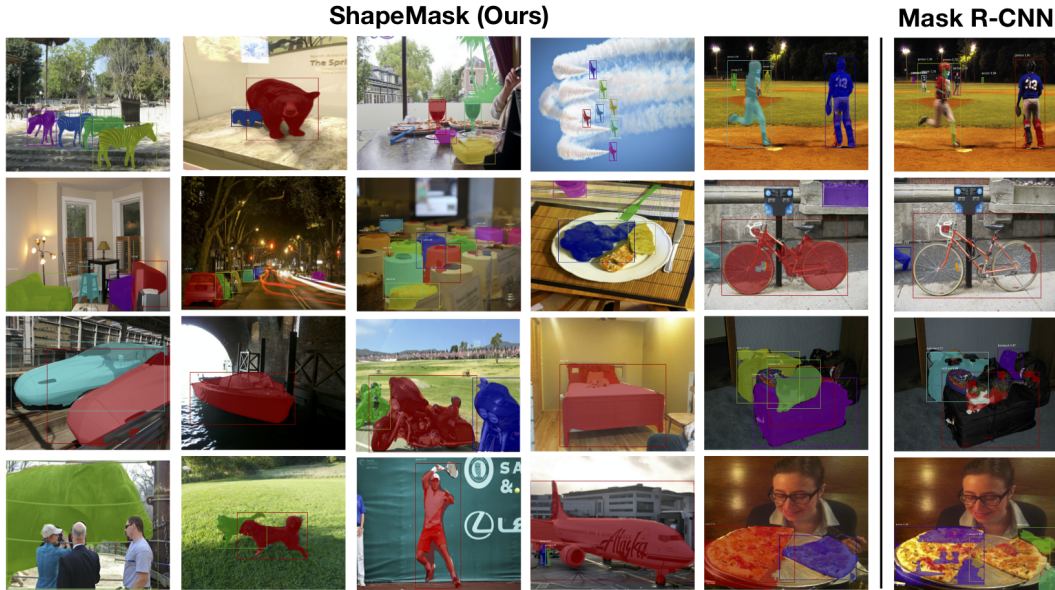


Figure 7: Visualization of ShapeMask on novel categories. For clarity, we only visualize the masks of novel categories. ShapeMask is able to segment many challenging objects well without seeing mask annotations in the same categories. It learns to predict object-like shapes for novel categories in many cases where Mask R-CNN does not (see rightmost column).

backbone	method	voc $\rightarrow$ non-voc						non-voc $\rightarrow$ voc					
		AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
FPN	Mask R-CNN [21]	18.5	34.8	18.1	11.3	23.4	21.7	24.7	43.5	24.9	11.4	25.7	35.1
	Our Mask R-CNN	21.9	39.6	21.9	16.1	29.7	24.6	27.2	39.6	27.0	16.4	31.8	35.4
	GrabCut Mask R-CNN [21]	19.7	39.7	17.0	6.4	21.2	35.8	19.6	46.1	14.3	5.1	16.0	32.4
	Mask <sup>X</sup> R-CNN [21]	23.8	42.9	23.5	12.7	28.1	33.5	29.5	52.4	29.7	13.4	30.2	41.0
	Oracle Mask R-CNN [21]	34.4	55.2	36.3	15.5	39.0	52.6	39.1	64.5	41.4	16.3	38.1	55.1
	Our Oracle Mask R-CNN	34.3	54.7	36.3	18.6	39.1	47.9	38.5	64.4	40.4	18.9	39.4	51.4
FPN	ShapeMask (ours)	<b>30.2</b>	<b>49.3</b>	<b>31.5</b>	<b>16.1</b>	<b>38.2</b>	<b>38.4</b>	<b>33.3</b>	<b>56.9</b>	<b>34.3</b>	<b>17.1</b>	<b>38.1</b>	<b>45.4</b>
	Oracle ShapeMask (ours)	35.0	53.9	37.5	17.3	41.0	49.0	40.9	65.1	43.4	18.5	41.9	56.6
NAS-FPN [13]	ShapeMask (ours)	<b>33.2</b>	<b>53.1</b>	<b>35.0</b>	<b>18.3</b>	<b>40.2</b>	<b>43.3</b>	<b>35.7</b>	<b>60.3</b>	<b>36.6</b>	<b>18.3</b>	<b>40.5</b>	<b>47.3</b>
	Oracle ShapeMask (ours)	37.6	57.7	40.2	20.1	44.4	51.1	43.1	67.9	45.8	20.1	44.3	57.8

Table 1: Performance of ShapeMask (class-agnostic) on novel categories. At the top, voc  $\rightarrow$  non-voc means “train on masks in voc, test on masks in non-voc”, and vice versa. ShapeMask outperforms the state-of-the-art method Mask<sup>X</sup> R-CNN [21] by 6.4 AP on voc to non-voc transfer, and 3.8 AP on non-voc to voc transfer using the same ResNet backbone. ShapeMask has smaller gap with the oracle upper-bound than Mask<sup>X</sup> R-CNN. By using a stronger feature pyramid from [13], ShapeMask outperforms Mask<sup>X</sup> R-CNN by 9.4 and 6.2 AP.

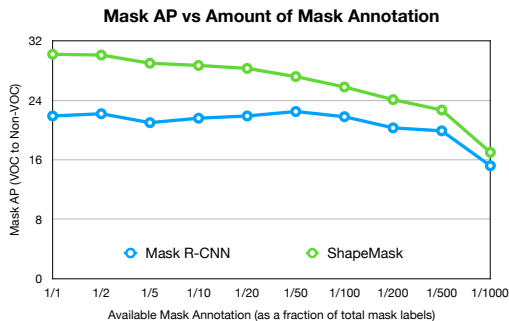


Figure 8: Generalization with less data. ShapeMask generalizes well down to 1/1000 of the training data.

## 4.2. Fully Supervised Instance Segmentation

Although the focus of ShapeMask is on generalization, this section shows that it is competitive as a general purpose instance segmentation algorithm.

**Main results:** We compare class-specific ShapeMask to leading instance segmentation methods on COCO in Table 2. Following previous work [19], training is on COCO train2017 and testing is on test-dev2017.

Using the same ResNet-101-FPN backbone, ShapeMask outperforms Mask R-CNN by 1.7 AP. With a stronger backbone, ShapeMask outperforms the best Mask R-CNN and MaskLab numbers by 2.9 and 2.7 AP. Since the focus of ShapeMask is to generalize to novel categories, we do not apply the techniques reported in [6, 34], including atrous

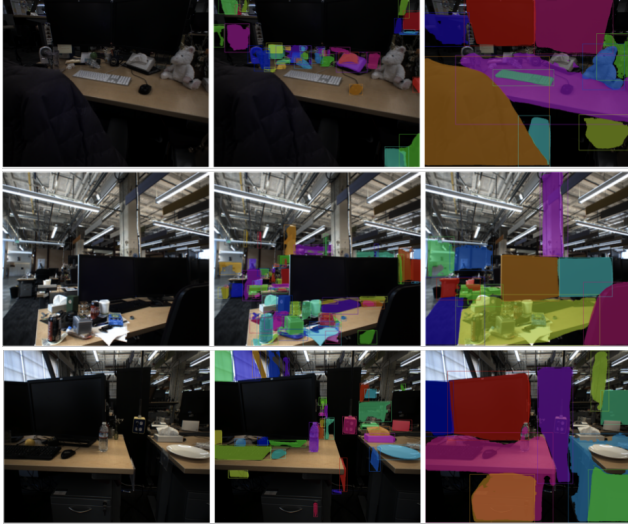


Figure 9: ShapeMask applied for object instance segmentation for robotics grasping. Here the ShapeMask model is trained on the COCO dataset and is not fine-tuned on data from this domain. As seen, it successfully segments the object instances, including novel objects such as a plush toy, a document, a tissue box, etc. For better visualization, smaller segmented objects are shown in the middle column and larger ones in the right column.

convolution, deformable crop and resize, mask refinement, adaptive feature pooling, heavier head, etc. Without any of these, ShapeMask ranks just behind PANet by 2.0 AP. Similarly, ShapeMask achieves 45.4 AP for box detection task without using the techniques reported by [5, 45, 34] – only second to the 47.4 AP of PANet (see Supp. Materials). Figure 1 of Supp. Materials visualizes the results of ShapeMask to demonstrate its ability to capture detailed contours, thin structures, and overlapping objects.

We benchmark the training and inference time with existing systems. Our training time of 11 hours on TPUs is 4x faster than all versions of Mask R-CNN [19, 15]<sup>2</sup>. For ResNet-101 model, we report competitive inference time among leading methods, where we note that our CPU time is unoptimized and can be reduced with more engineering. Among the heavier models, ShapeMask is the only method with reported runtimes. Training finishes within 25 hours on TPUs and runs at 5 fps per  $1024 \times 1024$  image on GPU. The Supp. Materials further show that by reducing the feature channels of mask branch, we can reduce the mask branch capacity by 130x and run 6x faster there (4.6ms) with marginal performance loss. These results show that ShapeMask is among the most efficient methods.

**Analysis of robust segmentation:** With pixelwise prediction approaches such as [19], the fate of mask is designed to depend heavily on detection quality. When detections are not reliable, there exists no mechanism for the mask branch

<sup>2</sup>[github.com/facebookresearch/Detectron/blob/master/MODEL\\_ZOO.md](https://github.com/facebookresearch/Detectron/blob/master/MODEL_ZOO.md)

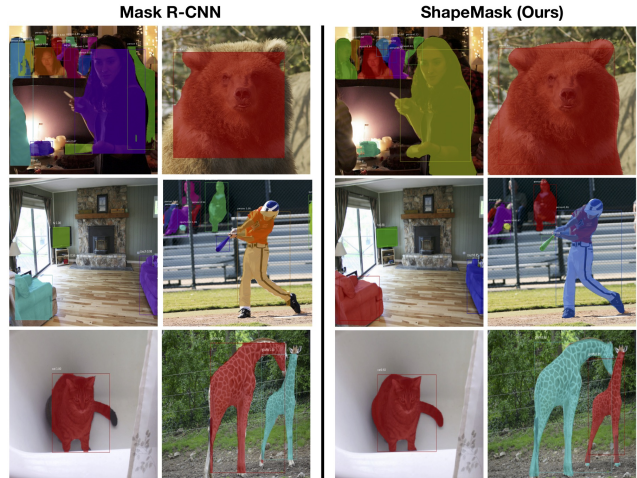


Figure 10: Analysis of Robust Segmentation. We stress-test Mask R-CNN and ShapeMask on randomly perturbed boxes (both were trained on whole boxes). Using soft detection priors, ShapeMask can handle poorly localized detections at test time while Mask R-CNN fails to do so by design of tight feature cropping.

to recover. In ShapeMask, masks are not confined to come from within detection boxes. We analyze the robustness of segmentation by conducting the following experiment.

First, we perturb the box detections at inference time by downsizing the width and height independently with a random factor  $x \sim U(0.75, 1.00)$ , where  $U$  represents uniform distribution. Downsizing avoids the complication of overlapping detections. Figure 10 compares the masks produced by Mask R-CNN and ShapeMask under this perturbation. Since Mask R-CNN can only produce masks within the boxes, it is not able to handle poorly localized detections. In contrast, ShapeMask uses detection merely as soft shape priors and manage to correct those cases without being trained for it at all. In addition, Table 3 quantifies the effect of downsized detections on mask quality. We see a significant drop in Mask R-CNN performance while ShapeMask remains stable. In addition, we show that training ShapeMask on downsized boxes improves its robustness.

### 4.3. Ablation Study

To understand our system further, we compare the uniform box prior with our learned detection prior, and the direct mask decoding [19] with our instance conditioned mask decoding. Table 4 shows our partially supervised system ablation results on COCO val2017 using ResNet-101-FPN. Surprisingly, using either object shape prior or instance embedding greatly improves from the baseline by about 12 and 5 AP. Combining both techniques boosts the performance even further. Similar results are found for the fully supervised setting (Supp. Materials).

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>	Training (hrs)	Inference (X + Y ms)	GPU
FCIS+++ [28] +OHEM	ResNet-101-C5-dilate	33.6	54.5	-	-	-	-	24	240	K40
Mask R-CNN [19]	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4	44	195 + 15	P100
Detectron Mask R-CNN [15]	ResNet-101-FPN	36.4	-	-	-	-	-	50	126 + 17	P100
ShapeMask (ours)	ResNet-101-FPN	37.4	58.1	40.0	16.1	40.1	53.8	11*	125 + 24	V100
Mask R-CNN [19]	ResNext-101-FPN	37.1	60.0	39.4	16.9	39.9	53.5	-	-	-
MaskLab [6]	Dilated ResNet-101	37.3	59.8	39.6	19.1	40.5	50.6	-	-	-
PANet [34]	ResNext-101-PANet	42.0	65.1	45.7	22.4	44.7	58.1	-	-	-
ShapeMask (ours)	ResNet-101-NAS-FPN [13]	40.0	61.5	43.0	18.3	43.0	57.1	25*	180 + 24	V100

Table 2: ShapeMask Instance Segmentation Performance on COCO. Using the same backbone, ShapeMask outperforms Mask R-CNN by 1.7 AP. With a larger backbone, ShapeMask outperforms Mask R-CNN and MaskLab by 2.9 and 2.7 AP respectively. Compared to PANet, ShapeMask is only 2.0 AP behind without using any techniques reported in [34, 6]. This shows that ShapeMask is competitive in the fully supervised setting. Timings reported on TPUs are marked with star signs. Inference time is reported following the Detectron format: X for GPU time, Y for CPU time. All mask APs are single-model, and are reported on COCO test-dev2017 without test time augmentation except Detectron on val2017 (gray).

Method	No Jittering	Jittering
Our Mask R-CNN	36.4	29.0
ShapeMask (ours)	37.2	<b>34.3</b>
ShapeMask w/ jittering training (ours)	37.2	<b>35.7</b>

Table 3: Instance segmentation Mask AP with jittered detections at test time. ShapeMask is more robust than Mask R-CNN by 5.3 AP. Adding jittering during training time makes ShapeMask more robust to it (last row).

Shape	Embed.	voc → non-voc			non-voc → voc		
		AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>
		13.7	28.0	12.0	24.8	45.6	23.5
	✓	26.2	44.6	27.1	29.4	51.7	29.0
	✓	26.4	44.9	27.2	30.6	53.4	30.4
	✓	30.2	49.3	31.5	33.3	56.9	34.3

Table 4: Ablation results for the partially supervised model.

#### 4.4. The Influence of Shape Priors

We conduct the following experiment to study how the quality of shape priors affects the final masks. We use the IoU of detection prior as a proxy for the distance to prior shapes in training set. This captures both the presence of similar shapes in the training set, and whether the shape priors are correctly predicted for downstream segmentation. We plot the detection prior IoU vs. the final mask IoU for non-voc classes with a model trained on voc categories in Figure 11 with visualization of various regimes. The plot shows clear positive correlation between the prior and final mask IoUs for the categories. We show IoU because it isolates the effect of mask prediction from object detection.

## 5. Conclusion

We introduce ShapeMask that uses shape priors and instance embeddings for better generalization to novel categories. ShapeMask significantly outperforms state-of-the-art in the cross categories setup. It is robust against inaccurate detections, competitive in the fully supervised setting, and runs efficiently for training and inference. We believe it

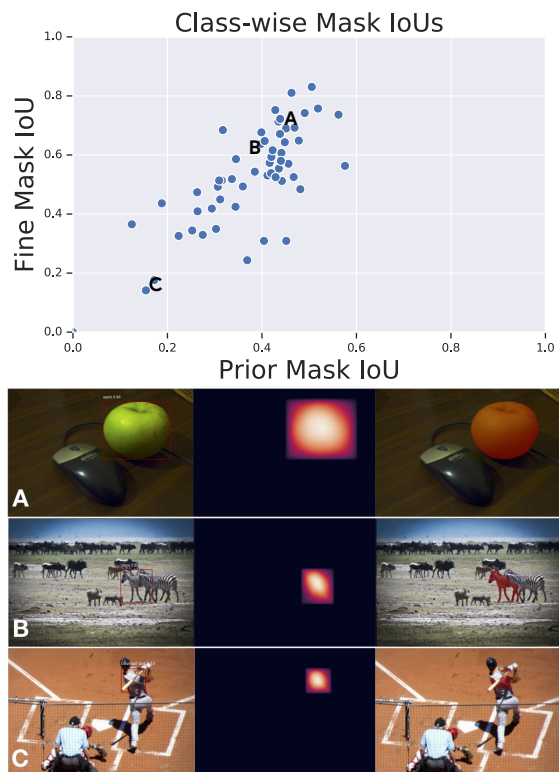


Figure 11: **Top:** Scatter plot of fine mask vs. detection prior mask IoU with the ground truth. Each dot represents a class average IoU. We observe a positive correlation among the classes. A, B, and C maps to the regime of very good, good, and poor mask IoUs. **Bottom:** Representative examples from regime A, B, and C. We observe that good priors tend to produce good masks (A and B), and a poor prior can cause the mask to go to the background (C).

is a step to further instance segmentation in the wild.

**Acknowledgements:** We want to thank Alexa Greenberg of X for engineering support to bring ShapeMask to robots, and Pengchong Jin for help with open-sourcing.



## References

- [1] Anurag Arnab and Philip HS Torr. Pixelwise instance segmentation with a dynamically instantiated network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 441–450, 2017. 1, 2
- [2] Min Bai and Raquel Urtasun. Deep watershed transform for instance segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2858–2866. IEEE, 2017. 1, 2
- [3] Eran Borenstein and Shimon Ullman. Learning to segment. In *ECCV*, 2004. 2
- [4] Bert De Brabandere, Davy Neven, and Luc Van Gool. Semantic instance segmentation with a discriminative loss function. In *Deep Learning for Robotic Vision Workshop at the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [5] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: delving into high quality object detection. *arXiv preprint arXiv:1712.00726*, 2017. 7
- [6] Liang-Chieh Chen, Alexander Hermans, George Papandreou, Florian Schroff, Peng Wang, and Hartwig Adam. Masklab: Instance segmentation by refining object detection with semantic and direction features. *arXiv preprint arXiv:1712.04837*, 2017. 2, 3, 5, 6, 8
- [7] Yi-Ting Chen, Xiaokai Liu, and Ming-Hsuan Yang. Multi-instance object segmentation with occlusion handling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3470–3478, 2015. 2
- [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 1, 2
- [9] Jifeng Dai, Kaiming He, Yi Li, Shaoqing Ren, and Jian Sun. Instance-sensitive fully convolutional networks. In *European Conference on Computer Vision*, pages 534–549. Springer, 2016. 5
- [10] Jifeng Dai, Kaiming He, and Jian Sun. Instance-aware semantic segmentation via multi-task network cascades. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3150–3158, 2016. 1, 2, 5
- [11] Achal Dave, Pavel Tokmakov, and Deva Ramanan. Towards segmenting everything that moves. *arXiv preprint arXiv:1902.03715*, 2019. 2
- [12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 5
- [13] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V. Le. NAS-FPN: learning scalable feature pyramid architecture for object detection. In *CVPR*, 2019. 5, 6, 8
- [14] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 5
- [15] Ross Girshick, Ilija Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. 7, 8
- [16] Google. *Cloud TPU*, 2019 (accessed March 12, 2019). <https://cloud.google.com/tpu/>. 2
- [17] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *European Conference on Computer Vision*, pages 297–312. Springer, 2014. 1, 2
- [18] Zeeshan Hayder, Xuming He, and Mathieu Salzmann. Boundary-aware instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5696–5704, 2017. 2
- [19] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017. 1, 2, 3, 4, 5, 6, 7, 8
- [20] Xuming He and Stephen Gould. An exemplar-based CRF for multi-instance object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 296–303, 2014. 2
- [21] Ronghang Hu, Piotr Dollár, Kaiming He, Trevor Darrell, and Ross Girshick. Learning to segment every thing. 1, 2, 3, 4, 5, 6
- [22] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre Luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, and Jonathan Ross. In-datacenter performance analysis of a tensor processing unit. 2017. 2
- [23] Anna Khoreva, Rodrigo Benenson, Jan Hendrik Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, volume 1, page 3, 2017. 1, 2, 3
- [24] Jaechul Kim and Kristen Grauman. Shape sharing for object segmentation. In *European Conference on Computer Vision*, pages 444–458. Springer, 2012. 2
- [25] Alexander Kirillov, Evgeny Levinkov, Bjoern Andres, Bogdan Savchynskyy, and Carsten Rother. Instancecut: from edges to instances with multicut. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5008–5017, 2017. 2
- [26] Shu Kong and Charless Fowlkes. Recurrent pixel embedding for instance grouping. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2
- [27] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalan-

- tidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017. 3
- [28] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. *arXiv preprint arXiv:1611.07709*, 2016. 2, 5, 8
- [29] Joseph J Lim, C Lawrence Zitnick, and Piotr Dollár. Sketch tokens: A learned mid-level representation for contour and object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3158–3165, 2013. 2
- [30] Tsung-Yi Lin, Piotr Dollár, Ross B Girshick, Kaiming He, Bharath Hariharan, and Serge J Belongie. Feature pyramid networks for object detection. In *CVPR*, volume 1, page 4, 2017. 5
- [31] Tsung-Yi Lin, Priyank Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 5
- [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2, 3, 5
- [33] Shu Liu, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *The IEEE International Conference on Computer Vision (ICCV)*, 2017. 1, 2
- [34] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8759–8768, 2018. 6, 7, 8
- [35] Yiding Liu, Siyu Yang, Bin Li, Wengang Zhou, Jizheng Xu, Houqiang Li, and Yan Lu. Affinity derivation and graph merge for instance segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 686–703, 2018. 1, 2
- [36] David Novotny, Samuel Albanie, Diane Larlus, and Andrea Vedaldi. Semi-convolutional operators for instance segmentation. 2018. 1
- [37] Trung Pham, Vijay BG Kumar, Thanh-Toan Do, Gustavo Carneiro, and Ian Reid. Bayesian semantic instance segmentation in open set world. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. 2
- [38] Pedro O Pinheiro, Ronan Collobert, and Piotr Dollár. Learning to segment object candidates. In *Advances in Neural Information Processing Systems*, pages 1990–1998, 2015. 1, 2
- [39] Pedro O Pinheiro, Tsung-Yi Lin, Ronan Collobert, and Piotr Dollár. Learning to refine object segments. In *European Conference on Computer Vision*, pages 75–91. Springer, 2016. 1, 2
- [40] Jordi Pont-Tuset, Pablo Arbelaez, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *IEEE transactions on pattern analysis and machine intelligence*, 39(1):128–140, 2017. 2
- [41] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 5
- [42] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *ACM transactions on graphics (TOG)*, volume 23, pages 309–314. ACM, 2004. 2
- [43] Alexander Sax, Bradley Emi, Amir R Zamir, Leonidas Guibas, Silvio Savarese, and Jitendra Malik. Mid-level visual representations improve generalization and sample efficiency for learning active tasks. *arXiv preprint arXiv:1812.11971*, 2018. 2
- [44] Xin Shu, Chang Liu, Tong Li, Chunkai Wang, and Cheng Chi. A self-supervised learning manipulator grasping approach based on instance segmentation. *IEEE Access*, 6:65055–65064, 2018. 1
- [45] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection–snip. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3578–3587, 2018. 7
- [46] David Weiss and Ben Taskar. Scalpel: Segmentation cascades with localized priors and efficient learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2035–2042, 2013. 2
- [47] Yi Yang, Sam Hallman, Deva Ramanan, and Charles C Fowlkes. Layered object models for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(9):1731–1743, 2012. 2
- [48] Yanzhao Zhou, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. Weakly supervised instance segmentation using class peak response. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3791–3800, 2018. 2, 3