

# StuffNet: Using ‘Stuff’ to Improve Object Detection

Samarth Brahmabhatt  
Georgia Institute of Technology  
samarth. robo@gatech.edu

Henrik I. Christensen  
UC San Diego  
hichristensen@ucsd.edu

James Hays  
Georgia Institute of Technology  
hays@gatech.edu

## Abstract

We propose a Convolutional Neural Network (CNN) based algorithm – StuffNet – for object detection. In addition to the standard convolutional features trained for region proposal and object detection [33], StuffNet uses convolutional features trained for segmentation of objects and ‘stuff’ (amorphous categories such as ground and water). Through experiments on Pascal VOC 2010, we show the importance of features learnt from stuff segmentation for improving object detection performance. StuffNet improves performance from 18.8% mAP to 23.9% mAP for small objects. We also devise a method to train StuffNet on datasets that do not have stuff segmentation labels. Through experiments on Pascal VOC 2007 and 2012, we demonstrate the effectiveness of this method and show that StuffNet also significantly improves object detection performance on such datasets.

## 1. Introduction

Recent advances in CNN-based object detection have resulted in significant accuracy increases. Faster R-CNN with its integrated region proposal network and end-to-end training scheme for both the object detection and region proposal networks sets an impressive baseline for many object detection datasets. In search for further improvement, the community has mainly chosen three paths: 1) *Deeper networks*: Most notably, ResNet [13] uses 101 layer deep CNNs trained using residual learning to learn more complex representations from images, 2) *Multi-task learning and multi-feature classification* (Figure 2): Multi-task learning refers to using different tasks to provide extra supervision and regularization for object detection. This is done by training the network for an auxiliary task in addition to object detection. The most notable examples of these auxiliary tasks are region proposal [10, 33], object semantic segmentation [9], object attribute prediction [8, 27] and face landmark prediction [39, 38] for face detection. Multi-feature classification refers to combining diverse features with features trained for object detection at the input of the

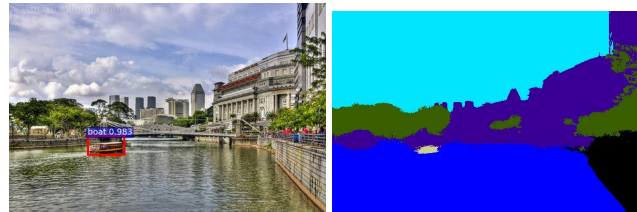


Figure 1: Most objects in images are surrounded by ‘stuff’ (e.g. the boat here is surrounded by water and building). We propose to allow stuff to influence the object detection decision. The figure shows the detection and semantic segmentation output of our model StuffNet-30 trained on Pascal VOC 2010 `trainval` on an image from MS COCO. See Figure 5e for legend.

classifier. These features capture complementary information like context [41], semantic segmentation of surrounding [28] and of objects [9]. Multi-task learning algorithms have multiple supervised outputs and typically increase performance for all tasks by learning shared features. In contrast, multi-feature learning algorithms have multiple complementary intermediate features, but don’t necessarily have extra supervision. 3) *Post-processing*: These approaches mainly try to better localize the objects by developing iterative localization algorithms that operate on object detections predicted by a base network and often achieve large performance gains [9].

This paper combines multi-task learning and multi-feature classification through a task which has been previously omitted from CNN-based object detection frameworks: segmentation of ‘stuff’. Stuff [15] refers to object categories which do not have a fixed shape and hence are difficult to bound with boxes e.g. water, sky, wall etc. We examine this task for two reasons: 1) It has been shown both in computer vision [15, 28, 41, 3] and cognitive psychology research [29] that identifying the local surrounding of an object helps to better identify it. This is because object identity tends to be correlated with spatially surrounding stuff, e.g. boats in most 2D images have water below them and sky/land above them (Figure 1). Especially for small,

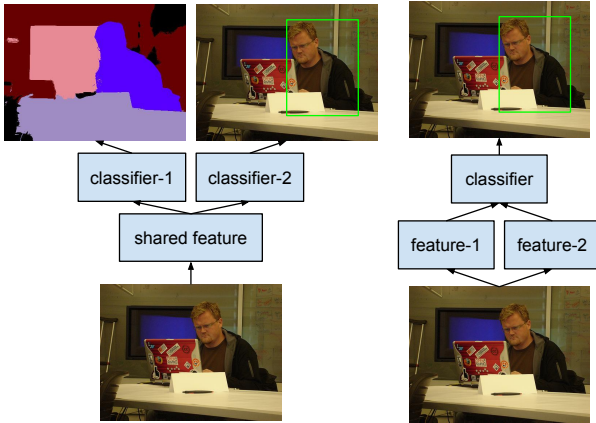


Figure 2: Multi-task learning (left) and multi-feature classification (right)

occluded, or unusually posed objects the local appearance might be ambiguous but context (in the form of surrounding stuff) can clarify what objects are likely. Hence, conditioning an object detection decision on features trained to identify stuff is likely to be beneficial. We demonstrate this through extensive evaluation in Section 4. Figure 4 also shows this qualitatively. 2) There have been many recent advancements in CNN-based semantic segmentation of images [25, 21, 24, 40]. Even though they are mainly focused on segmenting the objects, the network architectures and training mechanisms can be used just as easily for semantic segmentation of stuff. 3) Since the appearance of stuff (like sky, water, wall, trees, etc.) is similar across images from different datasets, a network trained to segment stuff using labels from one dataset can generalize well to a different dataset without re-training (we show this qualitatively in Figure 5). This means that once the stuff segmentation part of StuffNet is trained, we do not need stuff labels to train StuffNet on a different dataset. This is in contrast with most other tasks used for multi-task learning, and offsets the lack of large scale training data for stuff segmentation. Out of the numerous object detection datasets available to the community (Pascal VOC [5, 6, 7], MS COCO [22], ImageNet [19], Cityscapes [1], etc.), only Pascal 2010 (through Pascal Context [28]) and Cityscapes have stuff annotations. In this paper we train the stuff segmentation part of StuffNet only on the stuff labels in the Pascal Context dataset, and show that this can be deployed to Pascal VOC 2007 and Pascal VOC 2012 to boost object detection performance.

To summarize, we make the following contributions in this paper:

- *CNN architecture*: Our CNN architecture (Figure 3) supports both multi-task learning and multi-feature classification. We augment Faster R-CNN [33] with a semantic segmentation branch based on DeepLab [21].

This creates CNN that can learn both object detection and stuff segmentation in an end-to-end manner, allowing gradients from both tasks to influence a shared feature map (multi-task learning). Faster R-CNN uses convolutional feature maps to propose possible object regions and make decisions about the class of object present in those regions. These convolutional feature maps are trained solely for the tasks of region proposal and classification. In contrast, the StuffNet architecture conditions the object detection additionally on convolutional feature maps trained for stuff segmentation, and optionally, object segmentation (multi-feature classification). Through experiments on various datasets, we demonstrate that multi-task learning and multi-feature classification improves object detection performance.

- *Effective use of stuff labels*: Availability of stuff segmentation labels is a concern while training on a dataset. However, since stuff appears same across datasets, we only need one dataset with stuff labels to train the stuff segmentation part of StuffNet. We show that while training on a new dataset that does not have stuff labels, constraining the conv feature maps to produce the same segmentation output as a network trained on the earlier dataset acts as an effective regularizer and boosts object detection performance.

The rest of the paper is organized as follows: Section 2 reviews related work in this area. Section 3 describes our CNN architecture and training procedure. Section 4 shows results on object detection datasets and Section 5 discusses our results.

## 2. Related Work

### 2.1. CNN-based object detection

Fast R-CNN [10] pools convolutional feature maps using external region proposals (selective search [36]), and classifies them using fully connected layers interspersed with non-linearities. Faster R-CNN [33] incorporates the region proposal process into the network with a region proposal network. In this paper, we use Faster R-CNN as a baseline since it delivers state-of-the-art performance across many datasets (excluding deep residual networks) and has publicly available code. Recent advances with deep residual learning have resulted in 101 layer deep residual networks [13] that deliver superior performance. Even though we do not employ residual learning in this paper, our algorithm can be easily applied to residual nets to harness their deeper representations, since they have been applied to both object detection [13] and semantic segmentation [2].

## 2.2. Multi-task learning for improving object detection

Multi-task learning is the process of learning a representation that is useful for multiple tasks from the same input and has been found to improve the performance of both tasks [30, 37, 17, 31]. The most notable recent examples of multi task learning improving object detection are Fast R-CNN [10] and Faster R-CNN [33]. Fast R-CNN shows that training simultaneously for object detection and bounding box regression improves object detection performance. Faster R-CNN shows improvements in object detection performance by training for region proposal and bounding box regression. Recently, it has been shown that learning to predict facial landmarks along with learning to detect faces can improve face detection performance [39, 38].

## 2.3. Multi-feature classification for improving object detection

Recent works have shown that using features targeted at different tasks as input to the final classification layers of an object detector improves performance. For example, [28] uses counts of pixels around the object that are labeled with different stuff categories. [41] uses features from the penultimate layer of AlexNet extracted from a ring shaped region around the object. [9] uses convolutional features from a network trained separately for object segmentation, extracted from an expanded region proposal around the object. We compare the performance of StuffNet to [41] and [9] and show that StuffNet performs better.

## 3. StuffNet Architecture and Training

Stuffnet is a two-branch network which simultaneously performs object detection and semantic segmentation. The object detection branch is based on Faster R-CNN [33] and the semantic segmentation branch is based on DeepLab [21]. We first explain them briefly.

### 3.1. Faster R-CNN

**Why Faster R-CNN?** Faster R-CNN is the state of the art in object detection and integrates the region proposal and region classification tasks, which we traditionally performed separately. Most of the entries winning the MS COCO object detection challenge are based on Faster R-CNN, including [13], which uses a deep residual network to learn features but otherwise uses the same workflow. Hence Faster R-CNN makes a strong baseline to compare against. Moreover, the region proposal and classification modules of Faster R-CNN use the same input features. By training these features for diverse tasks (as described in Section 3.3), it is possible to influence both the region proposal and classification processes.

### 3.1.1 Faster R-CNN architecture

The network consists of a Region Proposal Network (RPN) [33] which proposes bounding boxes likely to contain an object, and the Fast R-CNN [10] detector that classifies the proposals (see bottom branch of Figure 3a). Both the RPN and Fast R-CNN share the `conv5-3` feature map from the VGG 16-layer CNN [35] as input.

The RPN is a fully convolutional network that produces a set of bounding boxes and their objectness scores. It slides a  $3 \times 3$  window over the `conv5-3` feature map and maps it from 1024 channels to 512 channels using a  $3 \times 3$  conv layer. Next, two  $1 \times 1$  conv layers are used in parallel. The `rpn-cls` layer predicts  $2k$  object presence/absence scores and the `rpn-reg` layer predicts  $4k$  co-ordinate regression parameters for  $k$  anchor proposal boxes. Each anchor is centered at the corresponding  $3 \times 3$  window, and has a unique scale and aspect ratio. The  $4k$  co-ordinate regression parameters are applied to the  $k$  anchors to produce the final region proposals. The  $2k$  object presence/absence scores are passed through a softmax layer to get the objectness scores.

The Fast R-CNN detector places an ROI pooling layer `roi-pool5` [14, 10] on top of `conv5-3`. This ROI pooling layer also takes the minibatch of region proposals as input. It adaptively max-pools the area of the activation map under each proposal into a grid equal to the input size of the following fully connected layer ( $7 \times 7$  for VGG-16 `fc6`). The final fully connected layer gives the object classification scores. Rectified activations of the `fc7` units are also passed to a bounding box regression layer `bbbox-pred`, which predicts parameters to change the four corners of the object proposal, separately for each class.

Faster R-CNN minimizes a weighted sum of four loss functions: 1) softmax loss for presence/absence of a foreground object in the region proposals, 2) smooth L1 loss [10] for anchor co-ordinate regression, 3) softmax loss for the classification of region proposals and 4) smooth L1 loss for region proposal co-ordinate regression.

### 3.1.2 Faster R-CNN training

We use the publicly available Python implementation of Faster R-CNN, which implements *joint approximate training*. The RPN training minibatch consists of 256 proposals from one image with a 1:1 ratio of positives (IoU with a ground-truth box  $> 0.7$ ) to negatives (IoU  $< 0.3$ ). The Fast R-CNN training minibatch has 128 proposals with a 1:3 ratio of foreground (IoU with a ground-truth box of the same class  $> 0.5$ ) to background ( $0.1 < \text{IoU}$  with any foreground ground-truth box  $< 0.5$ ). For more details on training parameters, please see [33].

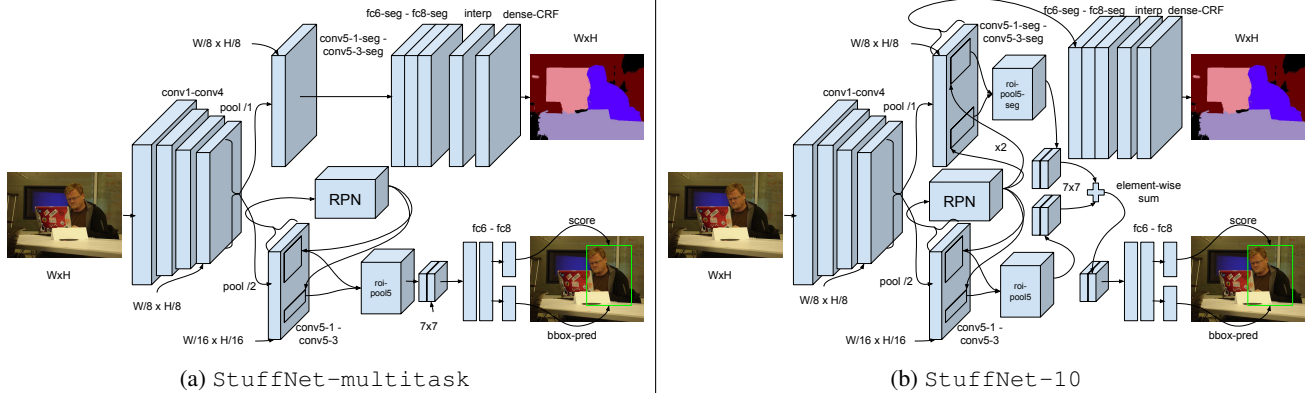


Figure 3: CNN architecture for two variants of StuffNet. StuffNet-multitask implements multi-task learning. StuffNet-10 implements both multi-task learning and use of diverse features. Please see [33] for details on the Region Proposal Network (RPN) and [21] for details on the dense CRF.

## 3.2. DeepLab

**Why DeepLab?** DeepLab [21] provides a fully convolutional [25] architecture based on the VGG 16-layer CNN that has been shown to produce state-of-the-art object segmentation results. It can be easily modified to train for stuff semantic segmentation by providing stuff segmentation labels for training and changing the channel dimension of the last  $1 \times 1$  convolution layer.

### 3.2.1 DeepLab architecture

In this paper, we use the Deeplab-LargeFOV [21] architecture (shown in the top branch of Figure 3a). DeepLab ‘convolutionizes’ the VGG 16-layer architecture [35] by using convolutions instead of fully connected layers  $fc6 - fc8$ . However, since spatial subsampling in max-pooling layers reduces the resolution of the feature map input to  $fc6$  by a factor of 32 for the VGG architecture, this leads to a very coarse segmentation output. DeepLab solves this problem by skipping spatial subsampling in max-pooling after the  $conv4-3$  and  $conv5-3$  layers and introducing ‘holes’ [26, 21] in the convolutional filters of layers downstream from  $conv4-3$ . Thus the resolution of the score map output by the DeepLab  $fc8$  layer is reduced by only a factor of 8. This score map is up-sampled using bilinear interpolation to get the final stuff segmentation output. At test time, DeepLab also applies a dense CRF [18] on its output to get crisp object boundaries. Finally, DeepLab minimizes a softmax loss over the segmentation label of each pixel in the upsampled score map.

### 3.2.2 DeepLab training

DeepLab uses a simple training procedure that involves separate training for the CNN and the dense CRF. The CNN is

trained using stochastic gradient descent with a minibatch of 10 images to provide unaries for the CNN. Subsequently, the parameters of the dense CRF are picked using cross validation, assuming the unary terms provided by the CNN are fixed. For more details, please see [21].

## 3.3. StuffNet

In this section we describe the changes we make to the architecture and training procedures of Faster R-CNN and DeepLab to construct StuffNet and train it in an end-to-end fashion.

### 3.3.1 StuffNet architecture

The objective of StuffNet is to improve object detection by utilizing 1) multi-task learning and 2) multi-feature classification. Towards this end, we first propose a simple combination of Faster R-CNN and DeepLab that implements multi-task learning, which we call StuffNet-multitask. Next we build on top of StuffNet-multitask and propose an architecture that incorporates semantic segmentation features into the detection decision process, which we call StuffNet-10 and StuffNet-30.

**Multi-task Learning.** As shown in Figure 3a, we combine Faster R-CNN and DeepLab to create StuffNet-multitask. Since both networks are based on the VGG 16-layer architecture, they share the conv layers up to  $conv4-3$ . This allows StuffNet-multitask to learn shared features that are useful for both object detection and semantic segmentation. Next, we pool the  $conv4-3$  feature map twice: once with a stride of 2 units ( $pool4-3$ ), and once with a stride of 1 unit ( $pool4-3-seg$ ) to produce the inputs to the next conv layers of Faster R-CNN and DeepLab respectively. This is essential to avoid a very

coarse resolution at the DeepLab output, as described in Section 3.2.1. The downstream layers of Faster R-CNN and DeepLab follow their respective architectures as described in Sections 3.1.1 and 3.2.1 respectively. Names of layers in the DeepLab branch have a `-seg` appended to distinguish them from corresponding layers in the Faster R-CNN branch. `StuffNet-multitask` minimizes the sum of the Faster R-CNN and DeepLab loss functions.

**Multi-feature classification.** `StuffNet-multitask` allows to learn a shared convolutional feature map that is used by both the object detection and semantic segmentation networks. However, it does not allow segmentation to directly influence the detection decision. We allow this influence in `StuffNet-10` and `StuffNet-30` by creating a connection between the two branches of `StuffNet` before the `fc6` layer, as shown in Figure 3b. We introduce another roi-pooling layer `roi-pool5-seg` that adaptively max-pools features from the `conv5-3-seg` layer. The regions input to `roi-pool5-seg` are double the size of those input to `roi-pool5`, to account for the higher resolution of the `conv5-3-seg` feature map. However, the output size of this roi-pooling layer is also set to  $7 \times 7$ , which allows us to combine the outputs of both roi-pooling layers by an *element-wise sum*. This combined representation is input to the `fc6` layer of the Faster R-CNN branch. We tried a non-linear combination of features as mentioned in [12], but found that a linear combination gives better performance and allows faster learning.

**Classes for semantic segmentation.** The main focus of this paper is to show that identifying the stuff around objects helps better identify them. Indeed, [28] have shown this using hand-crafted features. However, recently, some CNN-based object detection networks [9] have shown performance improvements by including features trained for the segmenting the same objects that are being detected. Hence we develop two variants of `StuffNet`: `StuffNet-10` segments the image into 10 stuff classes, while `StuffNet-30` segments the image into 10 stuff classes + 20 Pascal VOC object classes. We identify the 10 stuff classes from the main 33 classes used in the Pascal Context dataset [28]. They are: *wall, floor, water, tree, sky, road, ground, building, mountain* and *background*. Since the annotators for the Pascal Context dataset were free to create their own label names, we merge some categories into the 9 stuff categories mentioned above. Specifically, we merge *sidewalk* and *runway* into *road*, *ceiling* into *wall*, and *grass, platform, sand, snow* into *ground*.

### 3.3.2 StuffNet training

The learning rate policy, number of iterations and batch size differ greatly for Faster R-CNN [33] and DeepLab [21]. Faster R-CNN uses a batch size of 1 image (256 regions

for RPN, 128 regions for classification) and trains for 70k iterations. The learning rate starts at  $1e-3$  and follows ‘step’ learning rate policy, reducing by a factor of 10 after 50k iterations. On the other hand, DeepLab uses a batch size of 10 images, 20k iterations and a ‘poly’ learning rate with the power parameter set to 0.9 (see [21, 24] for details on the ‘poly’ learning rate policy). Since our main objective in this paper is to improve object detection performance, we use all the learning parameters from Faster R-CNN for all variants of `StuffNet`.

During training, `StuffNet` requires object bounding box labels as well as pixel-level semantic segmentation labels for all the segmentation classes (10 and 30 for `StuffNet-10` and `StuffNet-30` respectively). Hence, we first train `StuffNet-multitask`, `StuffNet-10` and `StuffNet-30` on the Pascal VOC 2010 dataset[6], which has object bounding box labels from the VOC object detection task, object segmentation labels from the VOC object segmentation task, and stuff labels from the Pascal Context dataset [28]. This network can then be used to generate ‘ground-truth’ stuff labels for datasets that lack them, as described next.

**Feature Constraining to train StuffNet without segmentation labels.** We use the fact that stuff appears same across datasets to develop a simple procedure to train `StuffNet` on datasets that lack stuff segmentation labels. We use the DeepLab branch of `StuffNet-10` trained on Pascal VOC 2010 to get *hallucinated* stuff ground truth labels for all images in that dataset. As shown qualitatively in Figure 5, the stuff labels produced in this manner are very close to ground truth, especially after the use of the dense CRF at the end of DeepLab to capture detailed region boundaries. We then train `StuffNet` as described previously, using these stuff labels. This constrains the conv layers shared between DeepLab and Faster R-CNN, and the conv layers in the DeepLab branch, to continue producing the same stuff labels as the training iterations proceed. We call this *feature constraining*. Figures 5c and 5d show the stuff labels produced by `StuffNet-10` on an image from MS COCO before and after this training respectively. It can be seen that this simple constraining strategy is effective in practice. An alternative to this method is to freeze the segmentation branch of the network and only re-train the detection branch. However, as shown in the multi-task setup of Figure 2, both branches operate on a shared representation. As this shared representation gets updated by the gradients from the detection branch, it no longer remains compatible with the frozen segmentation branch. We develop the feature constraining method to avoid this problem.

## 4. Results

We implemented `StuffNet` by modifying the publicly available Python code from Faster R-CNN, which uses the

Method	areo	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
Faster R-CNN [33]	79.9	78.4	66.6	44.7	<b>50.2</b>	77.9	70.3	86.0	40.6	60.4	50.1	<b>83.9</b>	76.7	<b>78.8</b>	77.4	36.0	68.8	55.4	75.6	64.5	66.1
StuffNet multitask	<b>80.3</b>	76.5	65.0	50.3	48.4	<b>78.3</b>	70.8	84.6	42.0	61.4	49.1	83.0	75.2	77.8	77.0	40.0	70.2	58.4	<b>79.5</b>	64.9	66.6
StuffNet-10	79.8	77.8	<b>67.7</b>	51.7	49.5	76.7	<b>72.1</b>	84.2	44.1	<b>64.5</b>	52.1	82.6	75.8	78.4	78.0	<b>42.2</b>	69.7	<b>58.9</b>	78.6	64.8	67.5
StuffNet-30	80.0	<b>79.5</b>	67.3	<b>53.5</b>	49.4	77.2	71.3	<b>86.4</b>	<b>46.4</b>	62.8	<b>54.4</b>	83.5	<b>78.8</b>	78.6	<b>78.5</b>	38.1	<b>70.8</b>	58.6	79.0	<b>65.4</b>	<b>68.0</b>

Table 1: Ablation study analysing effects of multi-task learning and feature diversification. Trained on Pascal VOC 2010 train, tested on Pascal VOC 2010 val

Method	Size	areo	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP	mAP change
Faster R-CNN [33]	small	34.0	<b>16.7</b>	7.8	27.6	<b>15.7</b>	<b>33.3</b>	42.8	0.0	0.0	29.6	0.0	0.0	25.0	<b>16.7</b>	30.2	<b>19.8</b>	43.3	0.0	<b>33.3</b>	0.0	18.8	0.0
StuffNet-10	small	<b>35.7</b>	0.0	<b>16.9</b>	<b>33.0</b>	13.4	<b>33.3</b>	<b>47.9</b>	<b>33.3</b>	<b>3.6</b>	46.3	0.0	<b>33.3</b>	<b>31.3</b>	<b>16.7</b>	<b>35.5</b>	19.4	<b>45.3</b>	0.0	<b>33.3</b>	0.0	<b>23.9</b>	<b>+5.1</b>
StuffNet-30	small	28.1	0.0	15.5	26.5	13.4	0.0	46.0	<b>33.3</b>	<b>3.6</b>	<b>61.5</b>	0.0	<b>33.3</b>	20.8	8.3	33.2	16.4	42.7	0.0	<b>33.3</b>	0.0	20.8	+2.0
Faster R-CNN [33]	medium	53.4	67.6	59.8	48.8	<b>50.1</b>	52.9	70.1	<b>55.5</b>	33.8	66.1	2.4	<b>74.5</b>	60.7	57.0	68.3	45.3	69.7	13.3	30.0	56.5	51.8	0.0
StuffNet-10	medium	<b>57.0</b>	65.1	62.6	54.3	48.8	54.4	73.3	52.2	34.5	<b>69.3</b>	<b>3.8</b>	68.0	59.7	<b>60.4</b>	<b>71.2</b>	<b>50.8</b>	71.0	<b>30.0</b>	25.2	56.3	53.4	+1.6
StuffNet-30	medium	56.8	<b>70.0</b>	<b>62.7</b>	<b>54.6</b>	49.7	<b>55.6</b>	<b>74.0</b>	51.1	<b>37.5</b>	63.4	1.2	69.9	<b>68.5</b>	58.0	70.9	48.8	<b>73.9</b>	19.4	<b>42.9</b>	<b>60.8</b>	<b>54.5</b>	<b>+2.7</b>
Faster R-CNN [33]	large	91.8	86.4	81.7	57.7	79.8	89.0	87.7	<b>91.0</b>	58.9	76.6	62.5	<b>90.5</b>	87.2	90.4	87.8	51.8	82.9	73.6	85.2	78.9	79.6	0.0
StuffNet-10	large	90.1	<b>87.6</b>	83.2	<b>65.4</b>	83.0	87.4	<b>89.2</b>	89.4	<b>65.1</b>	<b>82.0</b>	64.2	89.3	86.7	<b>91.1</b>	88.1	<b>63.0</b>	81.9	74.8	<b>90.0</b>	<b>81.7</b>	<b>81.7</b>	<b>+2.1</b>
StuffNet-30	large	<b>92.6</b>	85.9	<b>83.3</b>	60.9	<b>83.3</b>	<b>90.4</b>	89.1	90.7	64.0	76.8	<b>65.0</b>	89.1	<b>87.3</b>	87.8	<b>88.2</b>	57.1	<b>83.5</b>	<b>76.9</b>	89.6	78.4	81.0	1.4

Table 2: Analyzing StuffNet performance gains for objects of different sizes. Trained on Pascal VOC 2010 train, tested on Pascal VOC 2010 val

Dataset	Method	areo	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP	
VOC 07	YOLO [32]	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	63.4
	MR-CNN & S-CNN [9]	<b>76.8</b>	75.7	67.6	55.1	45.6	77.6	76.5	78.4	46.7	74.7	68.8	79.3	74.2	77.0	62.5	37.4	64.3	63.8	74.0	<b>74.7</b>	67.5	67.5
	SSD300 [23]	73.4	77.5	64.1	59.0	38.9	75.2	80.8	78.5	46.0	67.8	69.2	76.6	82.1	77.0	72.5	41.2	64.2	69.1	78.0	68.5	68.0	68.0
	SSD512 [23]	75.1	81.4	69.8	<b>60.8</b>	46.3	<b>82.6</b>	<b>84.7</b>	<b>84.1</b>	48.5	75.0	67.4	82.3	83.9	<b>79.4</b>	76.6	<b>44.9</b>	69.9	69.1	78.1	71.8	71.6	71.6
	OHEM [34]	71.2	78.3	69.2	57.9	46.5	81.8	79.1	83.2	47.9	76.2	68.9	83.2	80.8	75.8	72.7	39.9	67.5	66.2	75.6	75.9	69.9	69.9
	Faster R-CNN [33]	70.0	80.6	70.1	57.3	49.9	78.2	80.4	82.0	<b>52.2</b>	75.3	67.2	80.3	79.8	75.0	76.3	39.1	68.3	67.3	<b>81.1</b>	67.6	69.9	69.9
StuffNet-30	72.6	<b>81.7</b>	<b>70.6</b>	60.5	<b>53.0</b>	81.5	83.7	83.9	<b>52.2</b>	<b>78.9</b>	<b>70.7</b>	<b>85.0</b>	<b>85.7</b>	77.0	<b>78.7</b>	42.2	<b>73.6</b>	<b>69.2</b>	79.2	73.8	<b>72.7</b>	<b>72.7</b>	
VOC 10	segDeepM-16 layers [41]	82.3	75.2	67.1	50.7	49.8	71.1	69.5	<b>88.2</b>	42.5	71.2	50.0	85.7	76.6	81.8	69.3	41.5	71.9	<b>62.2</b>	73.2	<b>64.6</b>	67.2	67.2
	Faster R-CNN [33]	80.9	75.0	70.6	54.8	52.9	74.0	74.3	87.2	46.9	<b>72.6</b>	50.4	85.9	78.2	83.0	79.0	41.5	<b>74.1</b>	58.4	79.0	63.5	69.1	69.1
	StuffNet-10 †	82.3	75.9	<b>72.1</b>	<b>56.3</b>	54.0	<b>77.5</b>	<b>75.7</b>	86.9	47.9	72.1	52.2	<b>86.2</b>	78.4	83.0	<b>79.4</b>	<b>45.4</b>	72.4	58.5	78.2	64.3	69.9	69.9
	StuffNet-30 ‡	<b>83.0</b>	<b>77.1</b>	71.9	55.5	<b>54.5</b>	76.4	75.0	87.4	<b>48.7</b>	70.8	<b>54.6</b>	85.5	<b>79.3</b>	<b>83.2</b>	79.2	44.6	73.3	60.0	<b>79.8</b>	62.3	<b>70.1</b>	<b>70.1</b>
VOC 12	OHEM [34]	81.5	<b>78.9</b>	69.6	<b>52.3</b>	46.5	<b>77.4</b>	72.1	<b>88.2</b>	<b>48.8</b>	73.8	<b>58.3</b>	<b>86.9</b>	79.7	<b>81.4</b>	75.0	43.0	69.5	<b>64.8</b>	<b>77.5</b>	68.9	69.8	69.8
	Faster R-CNN [33]	82.3	76.4	71.0	48.4	45.2	72.1	72.3	87.3	42.2	73.7	50.0	86.8	78.7	78.4	77.4	34.5	70.1	57.1	71.1	58.9	67.0	67.0
	StuffNet-30 *	<b>83.0</b>	76.9	<b>71.2</b>	51.6	<b>50.1</b>	76.4	<b>75.7</b>	87.8	48.3	<b>74.8</b>	55.7	85.7	<b>81.2</b>	80.3	<b>79.5</b>	<b>44.2</b>	<b>71.8</b>	61.0	<b>78.5</b>	<b>65.4</b>	<b>70.0</b>	<b>70.0</b>

Table 3: StuffNet comparison with baselines on various datasets. For each dataset, the networks were trained on the respective trainval split and tested on the test split.

Method	wall	floor	water	tree	sky	road	ground	building	mountain	mAP
StuffNet-10	48.1	<b>39.4</b>	68.8	64.8	84.6	<b>39.0</b>	61.4	40.6	40.4	54.1
Separate	<b>48.6</b>	38.9	<b>69.3</b>	<b>64.9</b>	<b>85.1</b>	37.8	<b>62.3</b>	<b>41.3</b>	<b>40.6</b>	<b>54.3</b>
StuffNet-30	66.7	55.3	75.4	77.3	90.0	61.0	73.8	59.9	59.2	68.7

Table 4: Stuff segmentation performance results. Separate indicates a DeepLab network trained separately for segmenting the 10 stuff classes, using training parameters described in Section 3.3.2. StuffNet-10 and Separate were trained on Pascal VOC 2010 train and tested on Pascal VOC 2010 val. StuffNet-30 was trained on Pascal VOC 2012 trainval using the feature constraining method described in Section 3.3.2 and tested on Pascal VOC 2012 test.

Caffe library [16]. All our experiments use 70K training iterations. The learning rate starts at 1e-3 and drops by a factor of 10 after 50K iterations. We use a momentum of 0.9 and weight decay of 0.0005, as is standard practice. We initialize Faster R-CNN from the public VGG 16-layer model trained on ImageNet. Convolution layers specific to DeepLab are initialized using their public ‘init’ model, which is modified from VGG-16. The RPN and final fully connected layers are initialized randomly using

Xavier [11] initialization. In this section we evaluate the various StuffNet variants on their object detection performance on a variety of datasets. For all our experiments, we use Faster R-CNN as a common baseline.

#### 4.1. Impact of multi-task learning and multi-feature classification

StuffNet relies on 1) multi-task learning and 2) multi-feature classification to improve object detection performance. To analyse the impact of multi-task learning and multi-feature classification, we train the three StuffNet variants on Pascal VOC 2010 train split, and present evaluation results on the val split in Table 1. Multi-task learning (StuffNet-multitask) improves performance by 0.5% mAP, while its combination with multi-feature classification (StuffNet-10) improves performance by 1.5% mAP. Adding in object segmentation features (StuffNet-30) delivers an additional 0.5% mAP.

#### 4.2. Impact on different sized objects

Cognitive psychology [29] and computer vision [3] research has shown that identifying the surroundings is most helpful for identifying visually impoverished (i.e. blurred,

partially occluded or small in size) objects. To analyze the impact of multi-task learning and multi-feature classification on the detection performance of visually impoverished objects, we evaluated Faster R-CNN, *StuffNet-10* and *StuffNet-30*, considering objects of three different sizes. As in the MS COCO [22] dataset, the sizes for objects are defined as *small* (area < 32x32 sq. pixels), *medium* (32x32 < area < 96x96 sq. pixels) and *large* (area > 96x96 sq. pixels). When benchmarking on objects of each size, the ground truth labels for other sizes were marked ‘ignore’ while evaluating each size using the metric presented in [4]. Table 2 shows that the increase in mAP by *StuffNet* against Faster R-CNN is highest for *small* objects (5.1 percentage points mAP increase). This shows empirically that multi-task learning and multi-feature classification improves the object detection performance especially for small-sized objects, which are usually visually impoverished and need extra contextual information to be identified.

### 4.3. Evaluation on Datasets without stuff labels

One of the attractive features of *StuffNet* is that once the segmentation branch has been trained, it can be used to train *StuffNet* on any dataset that does not have stuff labels. This is done using the feature constraining mechanism described in Section 3.3.2. In this section, we train *StuffNet-30* using this method on Pascal VOC 2007 [5] and Pascal VOC 2012 [7] *trainval* splits, and test them on the respective *test* splits. Results presented in Table 3 show significant improvements delivered by *StuffNet-30* compared to Faster R-CNN. We also compare the *StuffNet-30* performance against two recent baselines that use contextual information and diverse features for object detection: Multi-region and segmentation-aware CNN (MR-CNN & S-CNN) from [9] and *segDeepM* [41]. MR-CNN & S-CNN combines features pooled from various contextual regions around the region proposal and features from a network trained *separately* to segment the same objects that are being detected. [9] also develops an iterative object re-localization strategy, that is applied as a post-processing step. Since object re-localization is beyond the scope of this paper, we compare *StuffNet* against the results from [9] *without* object re-localization. *segDeepM* diversifies the features input to the region classifier using features pooled from an expanded area around the region proposal, and extracted from the AlexNet [20] CNN trained for the ImageNet image recognition challenge.

As shown in Table 3, compared to Faster R-CNN *StuffNet-30* achieves a *performance gain of 2.8 percentage points on VOC 2007 and 3.0 percentage points on VOC 2012*, both of which have no stuff labels. In contrast, the comparable models from [9] and [41] perform worse than Faster R-CNN by around 2% mAP. This demonstrates that multi-task learning and constraining features to

output segmentations learnt from VOC 2010 act as an effective regularizer for the the object detection task. Figure 4 shows some examples of *StuffNet-30* detections on images from Pascal VOC 2010 *test*, compared to Faster R-CNN. Figure 5 shows detection (Faster R-CNN and *StuffNet-30*) and segmentation (*StuffNet-30*) results on images from the MS COCO dataset, using models trained on Pascal VOC 2010 *trainval*. Note that *StuffNet-30* is particularly better than Faster R-CNN at detecting small objects.

### 4.4. Semantic segmentation evaluation

Table 4 shows the performance of *StuffNet-10* and a separately trained DeepLab network on the task of segmenting 10 stuff classes. *StuffNet* performance is lower than separately trained DeepLab by 0.2 percentage points. We hypothesize that an appropriately weighted sum of the Faster R-CNN and DeepLab cost functions while training *StuffNet* will improve semantic segmentation performance. This is because both losses are normalized by the number of data-points contributing to the sum (number of region proposals in mini-batch for Faster R-CNN and number of pixels in the image for DeepLab). This normalization weighs the loss for a single pixel’s wrong semantic segmentation output with a very low factor. A careful search for the relative weights of the Faster R-CNN and DeepLab loss functions can be done using cross validation.

The last line of Table 4 shows the performance of a *StuffNet-30* trained with the feature constraining method on Pascal VOC 2012, a dataset that does not have stuff labels. The ground-truth labels for this evaluation were generated using *StuffNet-30* trained on Pascal VOC 2010 *trainval*. The high performance demonstrates that the feature constraining method is effective. Figure 4 shows some examples of object and stuff semantic segmentation using *StuffNet* models trained both normally and using feature constraining.

## 5. Conclusion

We presented *StuffNet*, an algorithm that uses multi-task learning and multi-feature classification to improve object detection, especially for small objects. *StuffNet* allows stuff and object semantic segmentation to closely influence the object detection process. The applicability of *StuffNet* can be extended to datasets that lack stuff or object segmentation labels by using a simple feature constraining procedure.

This paper opens up a few avenues for future research: 1) in the current architecture, the RPN does not benefit from multi-feature inputs. We conjecture that allowing this will improve localization performance, especially for small objects. 2) using deeper representations (e.g. ResNet [13]) should also further improve performance.

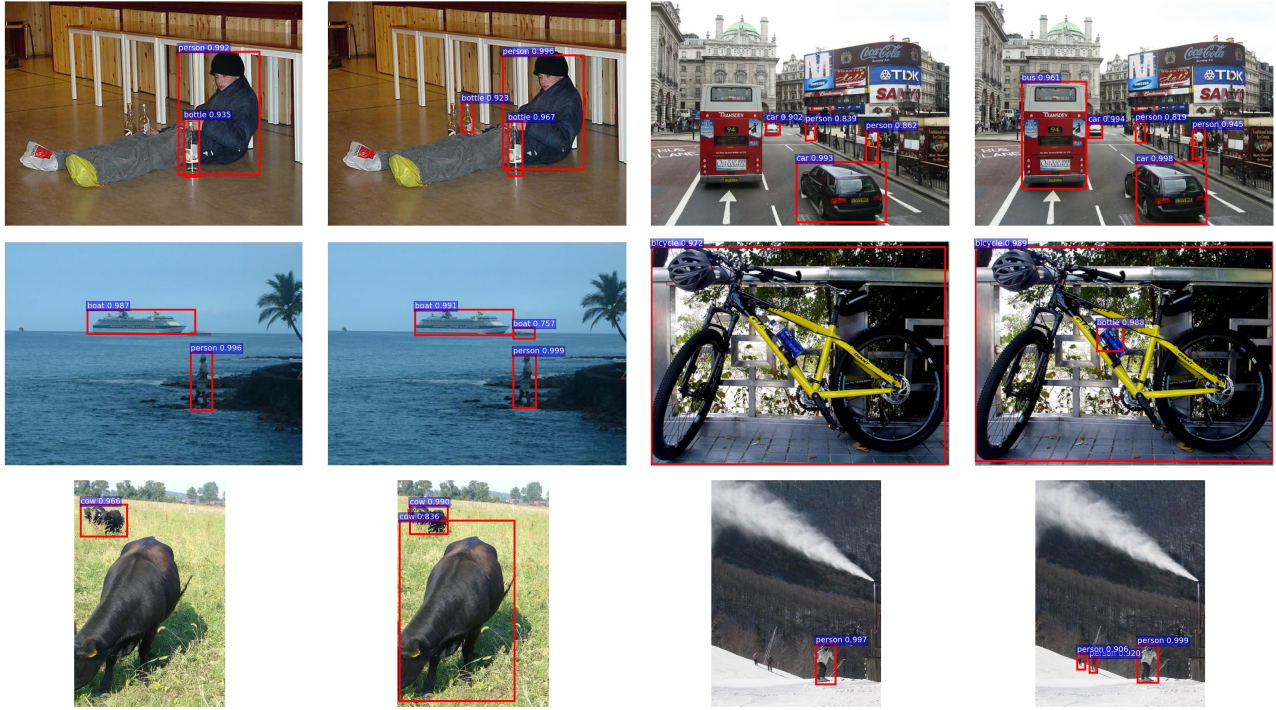


Figure 4: Pairs of Faster R-CNN (left) vs. StuffNet-30 (right) detections on sample images from the Pascal VOC 2010 test. Both models have been trained on Pascal VOC 2010 `trainval`

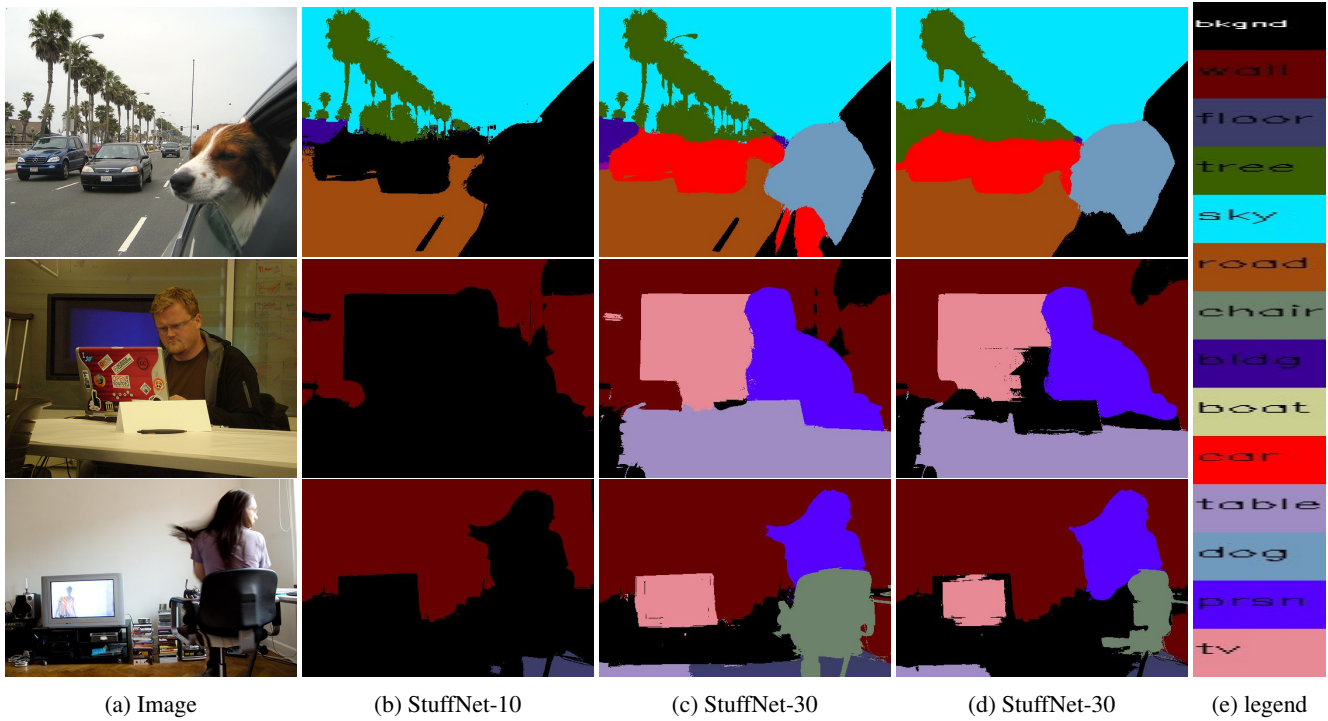
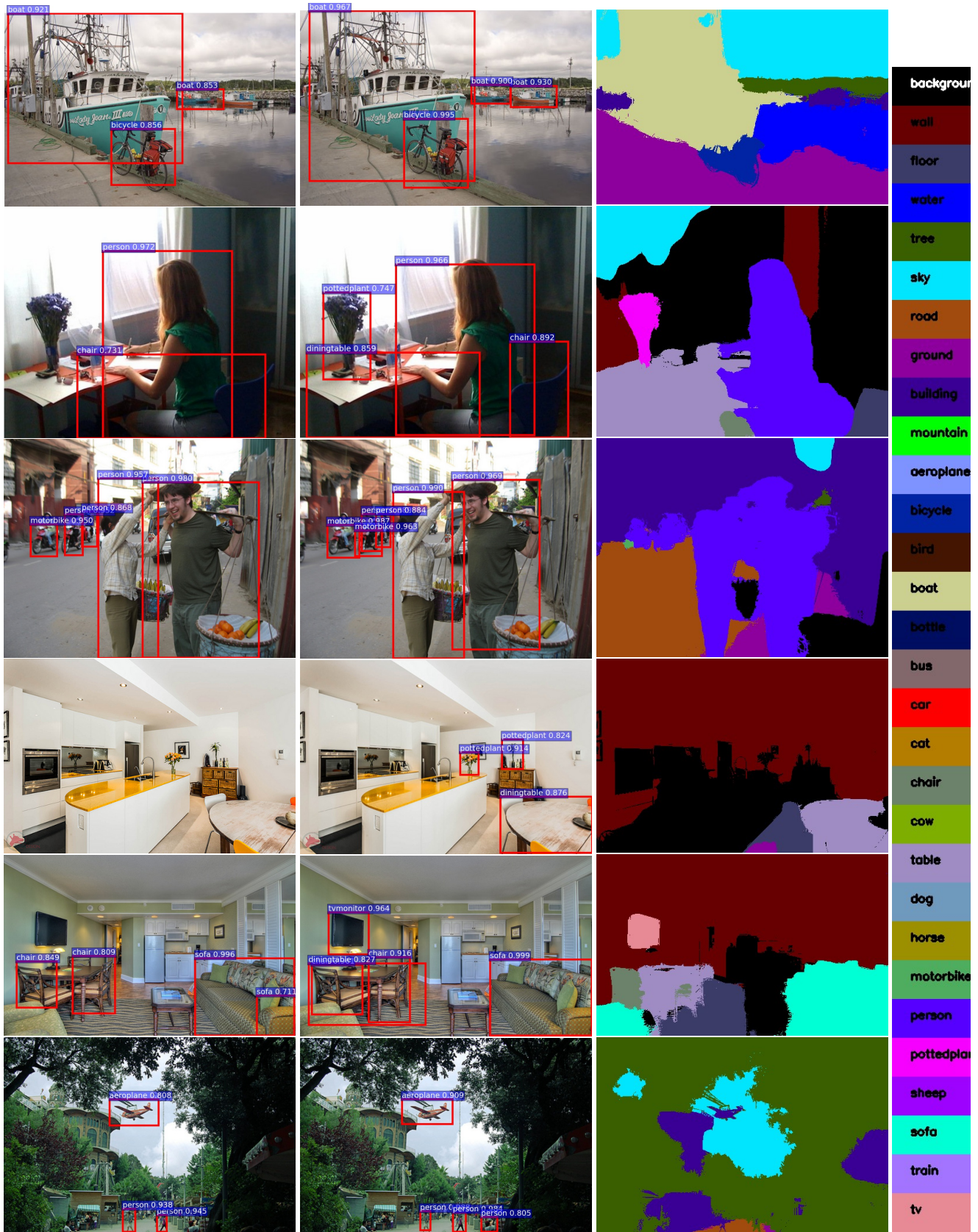


Figure 5: StuffNet segmentation qualitative results on images from the MS COCO dataset. Models (b) and (c) were trained on Pascal VOC 2010 `trainval`, while (d) was trained on Pascal VOC 2012 `trainval` using feature constraining.





(a) Faster R-CNN

(b) StuffNet-30

(c) StuffNet-30

(d) Legend

Figure 6: StuffNet-30 detection and segmentation qualitative results on images from the MS COCO dataset. The model was trained on Pascal VOC 2010 trainval. We also show Faster R-CNN detections in the first column for comparison. StuffNet-30 performs particularly well in detecting small objects.

## References

- [1] 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 *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [2] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. *arXiv preprint arXiv:1512.04412*, 2015. 2
- [3] S. K. Divvala, D. Hoiem, J. H. Hays, A. A. Efros, and M. Hebert. An empirical study of context in object detection. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1271–1278. IEEE, 2009. 1, 6
- [4] P. Dollar, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE transactions on pattern analysis and machine intelligence*, 34(4):743–761, 2012. 7
- [5] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>. 2, 7
- [6] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>. 2, 5
- [7] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. 2, 7
- [8] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 1778–1785. IEEE, 2009. 1
- [9] S. Gidaris and N. Komodakis. Object detection via a multi-region and semantic segmentation-aware cnn model. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1134–1142, 2015. 1, 3, 5, 6, 7
- [10] R. Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1440–1448, 2015. 1, 2, 3
- [11] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010. 6
- [12] S. Gupta, B. Hariharan, and J. Malik. Exploring person context and local scene context for object detection. *CoRR*, abs/1511.08177, 2015. 5
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 1, 2, 3, 7
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 37(9):1904–1916, 2015. 3
- [15] G. Heitz and D. Koller. Learning spatial context: Using stuff to find things. In *European conference on computer vision*, pages 30–43. Springer, 2008. 1
- [16] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 6
- [17] Z. Kang, K. Grauman, and F. Sha. Learning with whom to share in multi-task feature learning. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 521–528, 2011. 3
- [18] P. Krähenbühl and V. Koltun. Parameter learning and convergent inference for dense random fields. 4
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 2
- [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 7
- [21] C. Liang-Chieh, G. Papandreou, I. Kokkinos, K. Murphy, and A. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *International Conference on Learning Representations*, 2015. 2, 3, 4, 5
- [22] 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 *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 2, 7
- [23] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. *SSD: Single Shot MultiBox Detector*. Springer International Publishing, 2016. 6
- [24] W. Liu, A. Rabinovich, and A. C. Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015. 2, 5
- [25] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 2, 4
- [26] S. Mallat. *A wavelet tour of signal processing: the sparse way*. Academic press, 2008. 4
- [27] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert. Cross-stitch networks for multi-task learning. *arXiv preprint arXiv:1604.03539*, 2016. 1
- [28] 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. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 891–898, 2014. 1, 2, 3, 5
- [29] A. Oliva and A. Torralba. The role of context in object recognition. *Trends in cognitive sciences*, 11(12):520–527, 2007. 1, 6
- [30] A. Pentina, V. Sharmanska, and C. H. Lampert. Curriculum learning of multiple tasks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5492–5500, 2015. 3

- [31] A. Quattoni, M. Collins, and T. Darrell. Transfer learning for image classification with sparse prototype representations. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008. [3](#)
- [32] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [6](#)
- [33] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#)
- [34] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [6](#)
- [35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. [3](#), [4](#)
- [36] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013. [2](#)
- [37] J. Yim, H. Jung, B. Yoo, C. Choi, D. Park, and J. Kim. Rotating your face using multi-task deep neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 676–684, 2015. [3](#)
- [38] C. Zhang and Z. Zhang. Improving multiview face detection with multi-task deep convolutional neural networks. In *IEEE Winter Conference on Applications of Computer Vision*, pages 1036–1041. IEEE, 2014. [1](#), [3](#)
- [39] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In *European Conference on Computer Vision*, pages 94–108. Springer, 2014. [1](#), [3](#)
- [40] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1529–1537, 2015. [2](#)
- [41] Y. Zhu, R. Urtasun, R. Salakhutdinov, and S. Fidler. segdeepm: Exploiting segmentation and context in deep neural networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4703–4711, 2015. [1](#), [3](#), [6](#), [7](#)