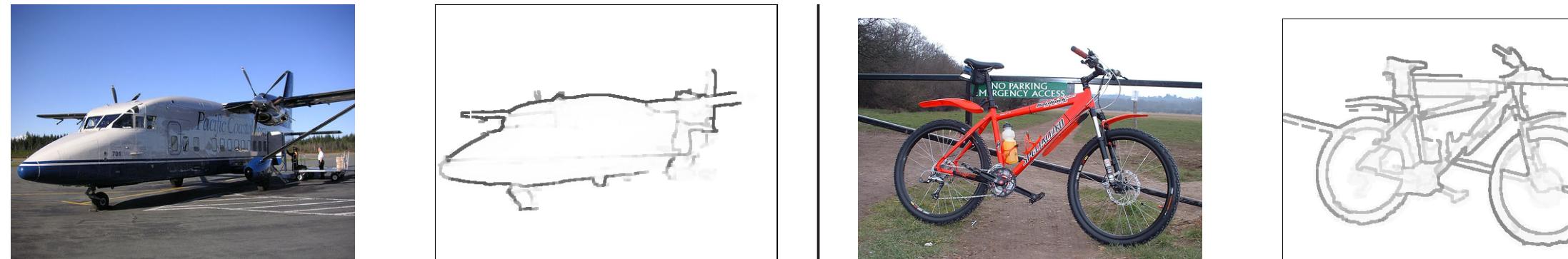


SEMANTIC BOUNDARY REFINEMENT BY JOINT INFERENCE FROM EDGES AND REGIONS

CHAO "HARRY" YANG UNIVERSITY OF SOUTHERN CALIFORNIA

DETECTING BOUNDARIES FOR SPECIFIC OBJECTS



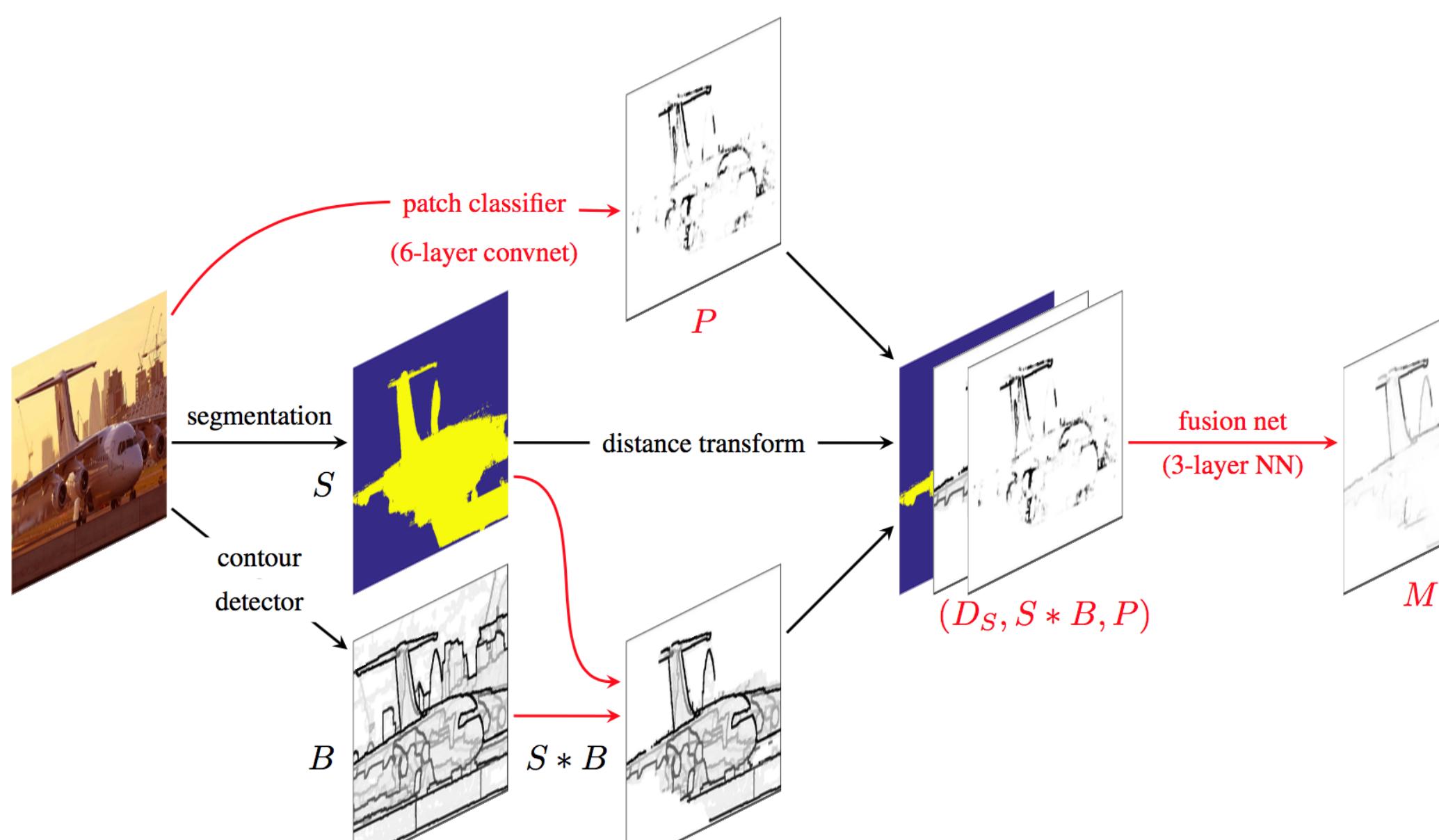
We study the problem of detecting boundaries for specific classes of objects. Our approach leverages recent advances in semantic segmentation and bottom-up boundary detection. We propose a mechanism for combining multiple sources of information: predicted segmentation masks, bottom-up contours, and a novel local class-specific boundary detector. These are jointly mapped to final category-specific boundary strength estimate by a trained classifier. In experiments on VOC2010 and Microsoft COCO data sets, our method dramatically outperforms recent prior work, for some classes doubling the accuracy of boundary prediction.

THE PIPELINE

Our system of category-specific boundary detection consists of the following components:

- Semantic segmentation (DeepLab).
- Contour detector (UCM).
- Local boundary detector (ConvNets)

A fusion component is used at the end to combine the outputs of the three components. The pipeline is illustrated in the following graph:



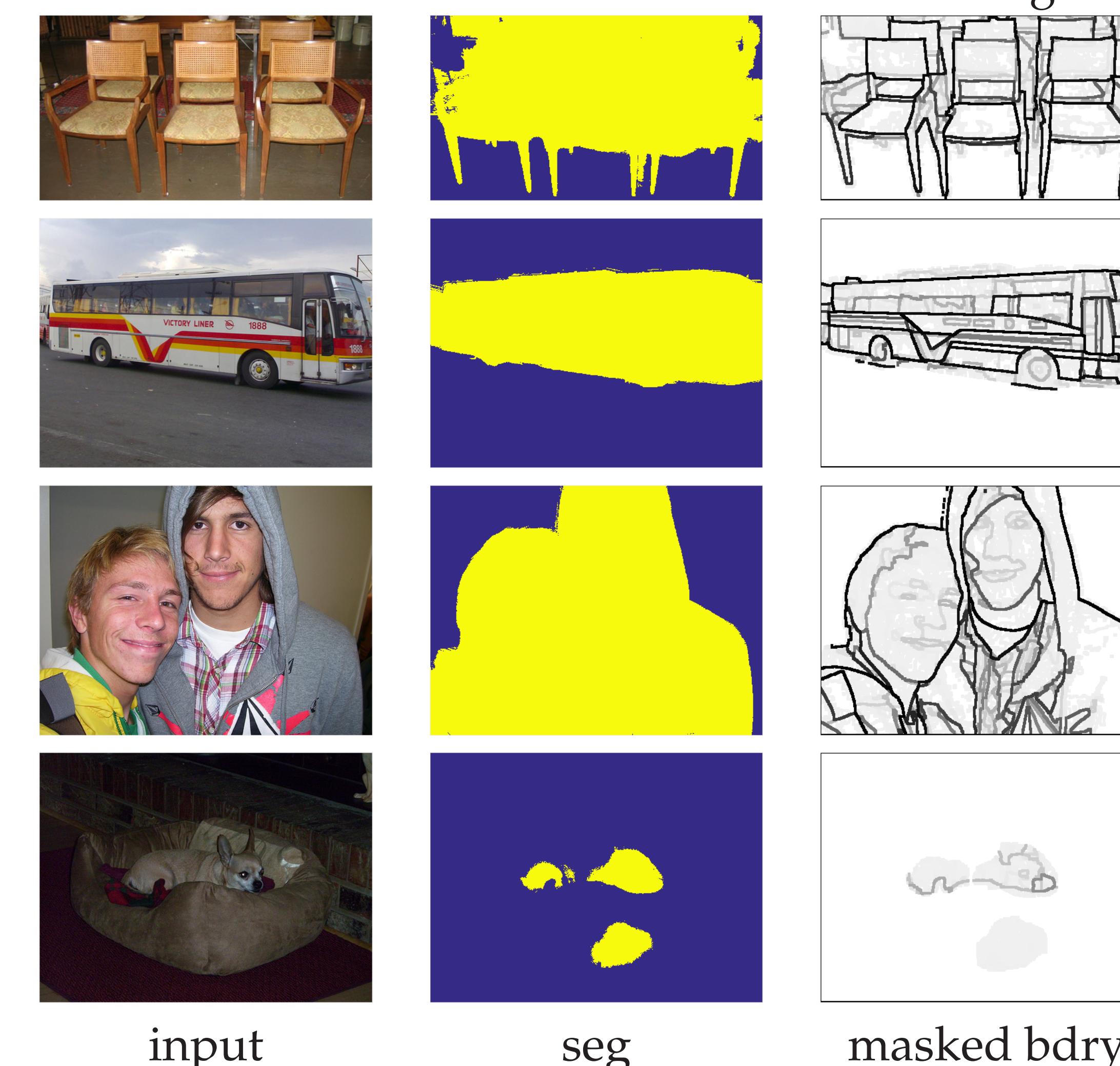
Our experiments use two datasets: VOC 2012 and COCO.

SEGMENTATION AND MASKED BOUNDARY BASELINE

Segmentation and masked boundary are the first two components of our system and are also used as the baseline.

- Segmentation: we use the semantic segmentation contour (DeepLab) as the boundary.
- Masked boundary: We intersect class agnostic boundary detection (UCM) with semantic segmentation mask (original segmentation dilated by 5 pixels).

The results are shown in the following figure:



LOCAL CLASS-SPECIFIC BOUNDARY DETECTOR

We can further improve the results by training a class-specific boundary predictor from scratch. We implement it as a multi-layer convnet, trained to classify image patches as boundary or non-boundary for the class at hand:

layer type	RF size	#units	stride
1 conv	4	96	1
2 max-pool	2		2
3 conv	3	256	1
4 max-pool	2		2
5 conv	4	64	1
7 ip		256	
8 ip		256	
9 ip		2	

INFORMATION FUSION

The three components described above each capture a different aspect of visual information.

- Segmentation: may offer poor localization of the boundaries.
- Masked boundary: low precision.
- local detector: based on limited view of the image due to its local nature.

We use a non-linear classifier to combine them:

$$\log \Pr(m = 1 | \mathbf{x}) \propto \mathbf{w}_3^T \sigma(\mathbf{w}_2^T \sigma(\mathbf{w}_1^T \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) + b_3,$$

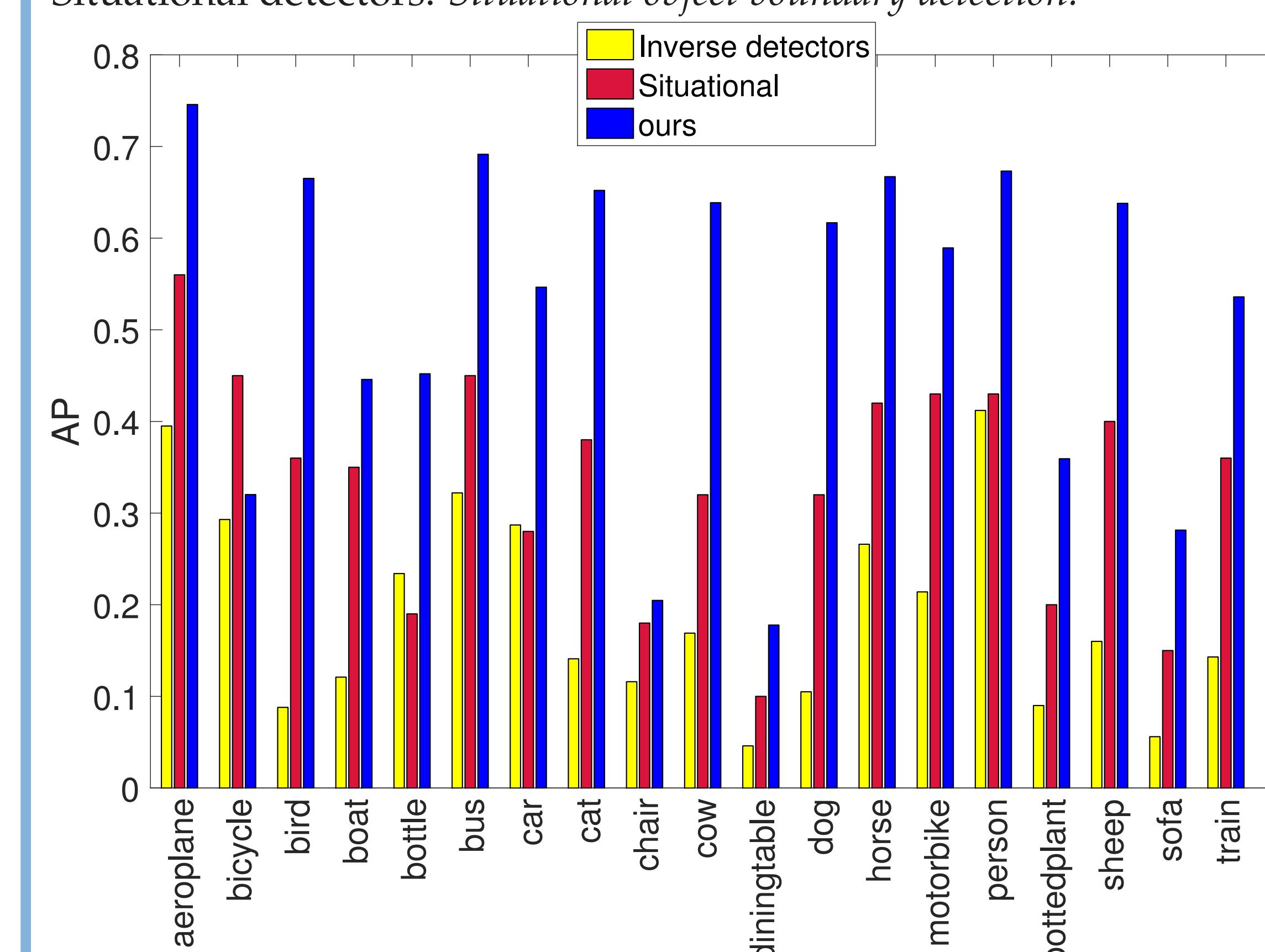
It can be viewed as a three-layer feedforward network with fully connected layers and sigmoid activation functions!

NUMERICAL RESULTS ON VOC2012

1. Per-class accuracy comparing with STOA:

Inverse detectors: *Semantic contours from inverse detectors.*

Situational detectors: *Situational object boundary detection.*

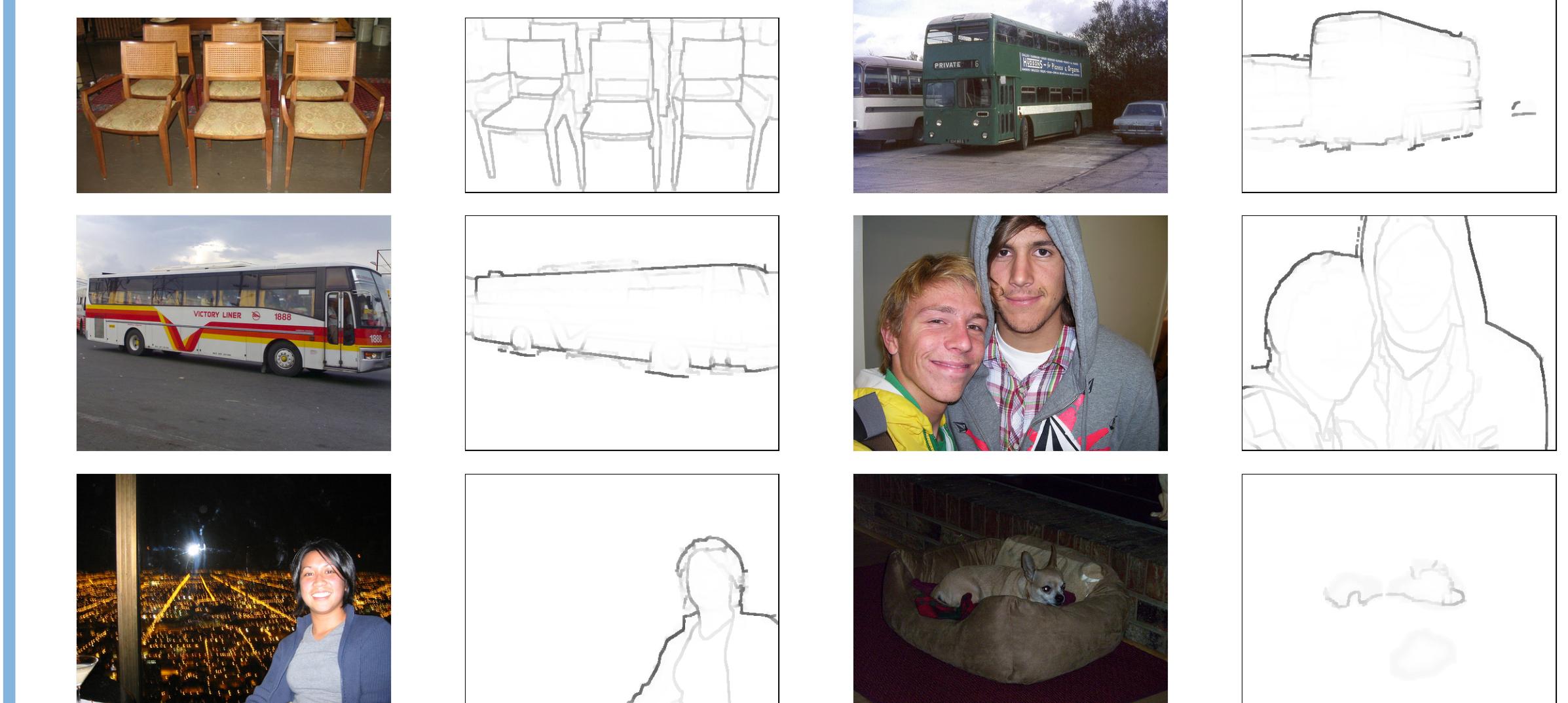


2. Mean average precision:

Method	Mean precision
Inverse detectors	19.9
Situational detectors	31.6
Segmentation only <i>S</i>	45.8
Masked boundary <i>S * B</i>	46.0
Local patch classifier <i>P</i>	26.9
Fusion (<i>D_s, S * B, P</i>)	51.8

We can see that our simple baseline already out-performs the inverse detectors and situational detectors by a large margin.

VISUAL RESULTS ON VOC2012



NUMERICAL RESULTS ON MICROSOFT COCO

we used the first 5,000 images from the val set as the test set for our experiments.

Method	Mean precision
Segmentation only <i>S</i>	36.9
Masked boundary <i>S * B</i>	38.5
Local patch classifier <i>P</i>	22.4
Fusion (<i>D_s, S * B, P</i>)	45.0

The relative standing reflects that in VOC results, but all the numbers are lower, reflecting the increased difficulty of COCO compared to VOC.

VISUAL RESULTS ON MICROSOFT COCO



CLASS AGNOSTIC RESULTS

We can also max-pooling the predictions of our category-specific detectors for each pixel, forming a class-agnostic boundary detection map:

