

## Summary

### Problem statement

- Object bounding boxes provide a definite background along with the extent of each object, but each box contains a mixture of foreground and background
- Goal: train segmentation models with object bounding boxes

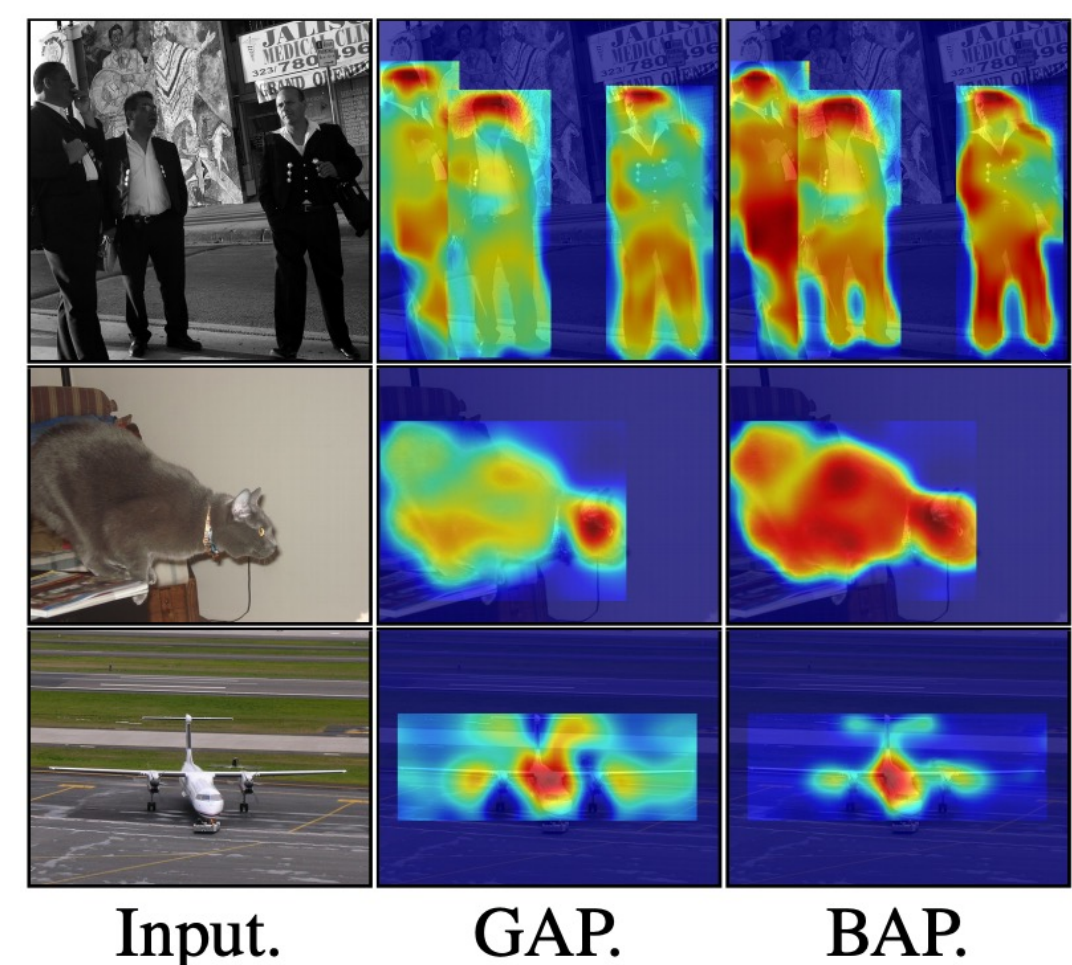
### Motivation

- How can we generate high-quality pseudo pixel-level labels from object bounding boxes?
- How can we train segmentation models with noisy labels?

### Contributions

- Introduce a simple yet effective framework which mainly consists of three stages: (1) train a CNN for image classification; (2) generate pseudo pixel-level labels; (3) train a CNN for semantic segmentation
- Propose a background-aware pooling (BAP) layer that leverages a background prior to separate foreground and background regions inside object bounding boxes
- Introduce a noise-aware loss (NAL) that alleviates the influence of incorrect labels adaptively
- Demonstrate state-of-the-art performance on PASCAL VOC 2012 and MS-COCO

## Comparison between GAP and BAP

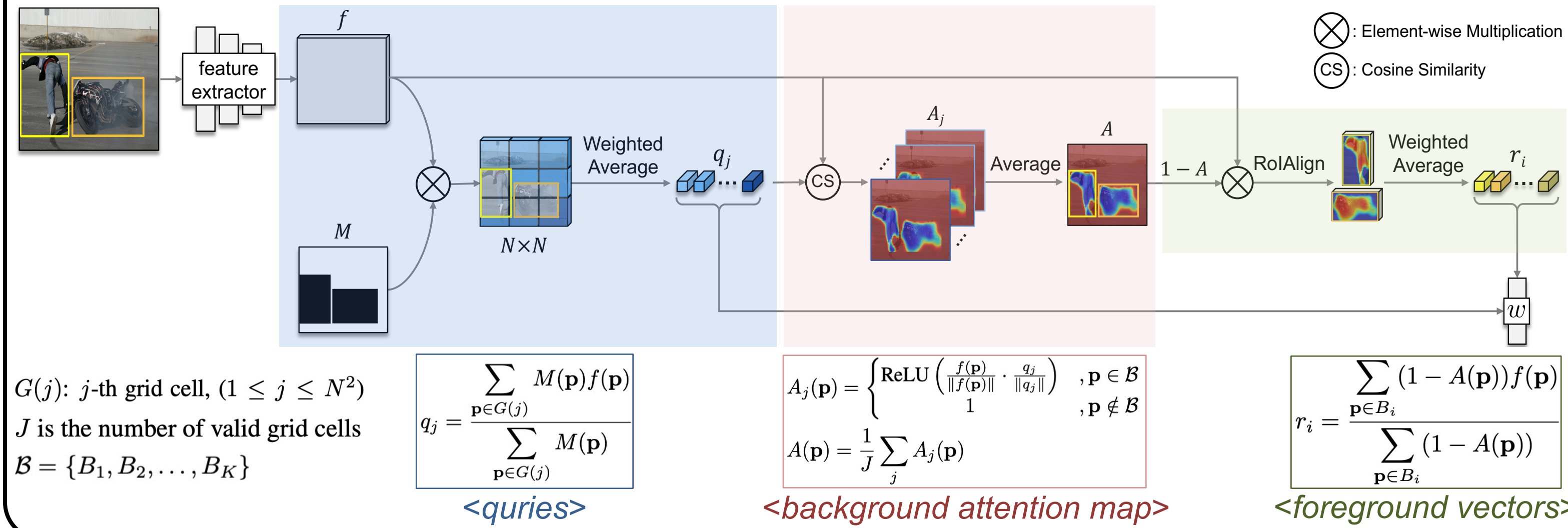


Method	train	val
<i>Supervision: Image-level labels (10K)</i>		
AffinityNet (CVPR 2018)	59.7	-
<i>Supervision: Boxes (10K)</i>		
Box	65.4	62.2
GrabCut (SIGGRAPH 2004)	65.7	66.1
MCG (PAMI 2016)	66.2	66.9
WSSL (ICCV 2015)	69.7	71.1
Ours		
GAP	75.5	76.1
BAP: $Y_{crf}$ w/o $u_0$	77.0	77.8
BAP: $Y_{crf}$	<b>78.7</b>	<b>79.2</b>
BAP: $Y_{ret}$	70.8	69.9

mIoUs of pseudo labels on PASCAL VOC 2012

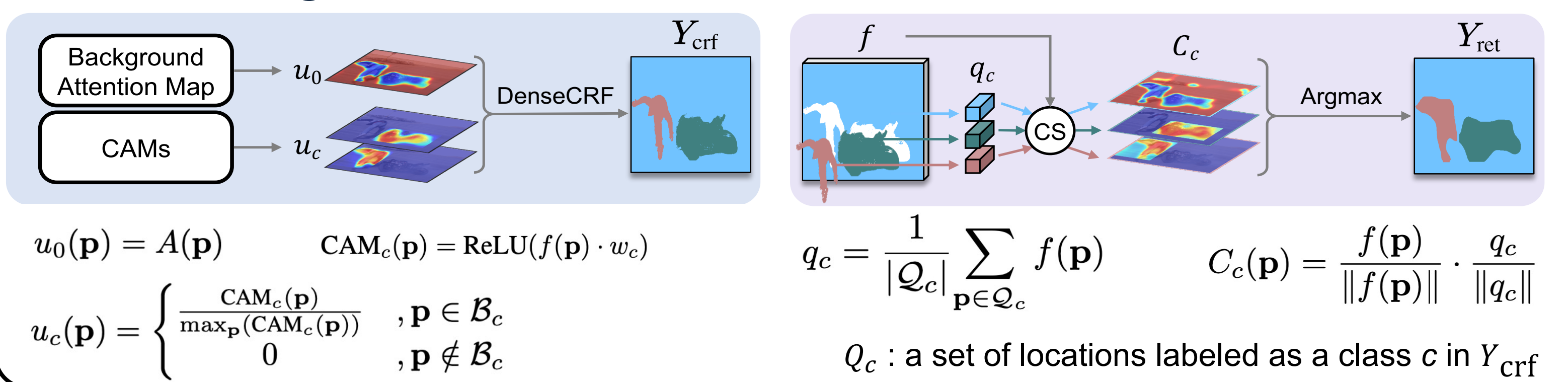
## Stage 1

### Image classification using BAP



## Stage 2

### Pseudo label generation



## Stage 3

The NAL is defined as follows:  $\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{wce}$

$$\mathcal{L}_{ce} = -\frac{1}{\sum_c |S_c|} \sum_c \sum_{\mathbf{p} \in S_c} \log H_c(\mathbf{p})$$

$$H_c(\mathbf{p}) = \frac{e^{\tau \frac{\phi(\mathbf{p})}{\|\phi(\mathbf{p})\|} \cdot \frac{W_c}{\|W_c\|}}}{\sum_i e^{\tau \frac{\phi(\mathbf{p})}{\|\phi(\mathbf{p})\|} \cdot \frac{W_i}{\|W_i\|}}$$

$H_c$ : a probability for a class  $c$

$S_c$ : a set of locations where both  $Y_{crf}$  and  $Y_{ret}$  give the same label  $c$

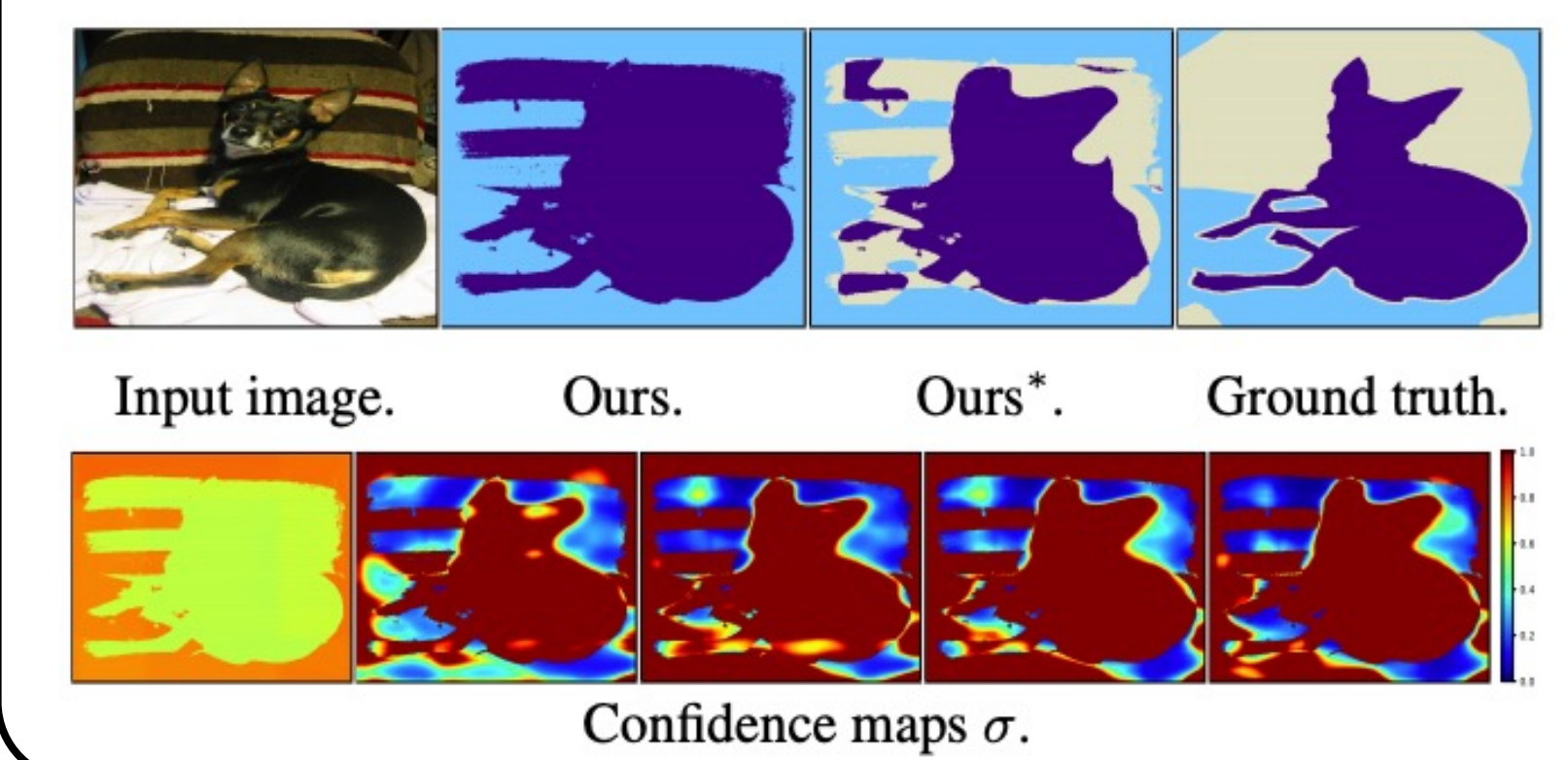
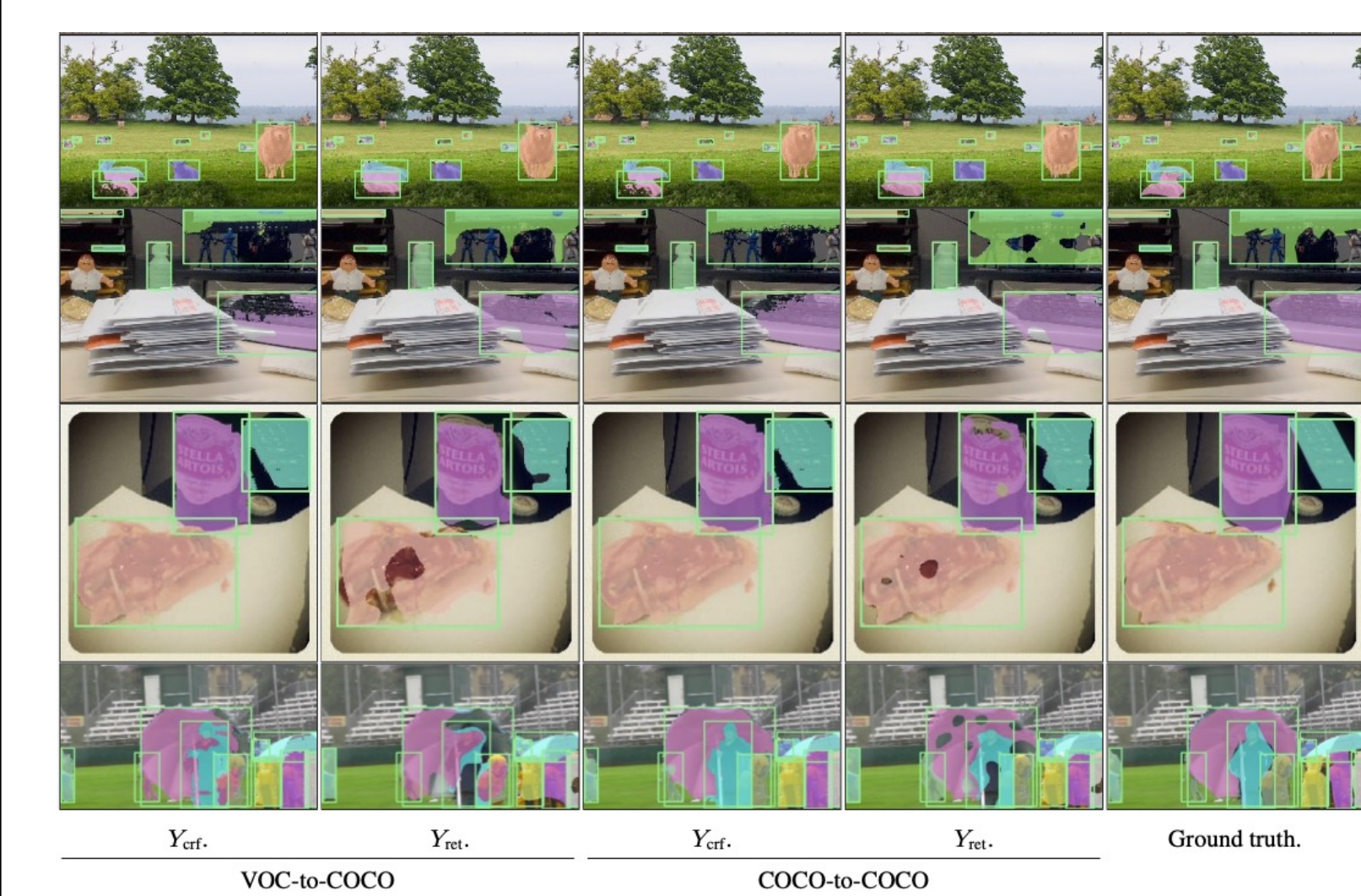
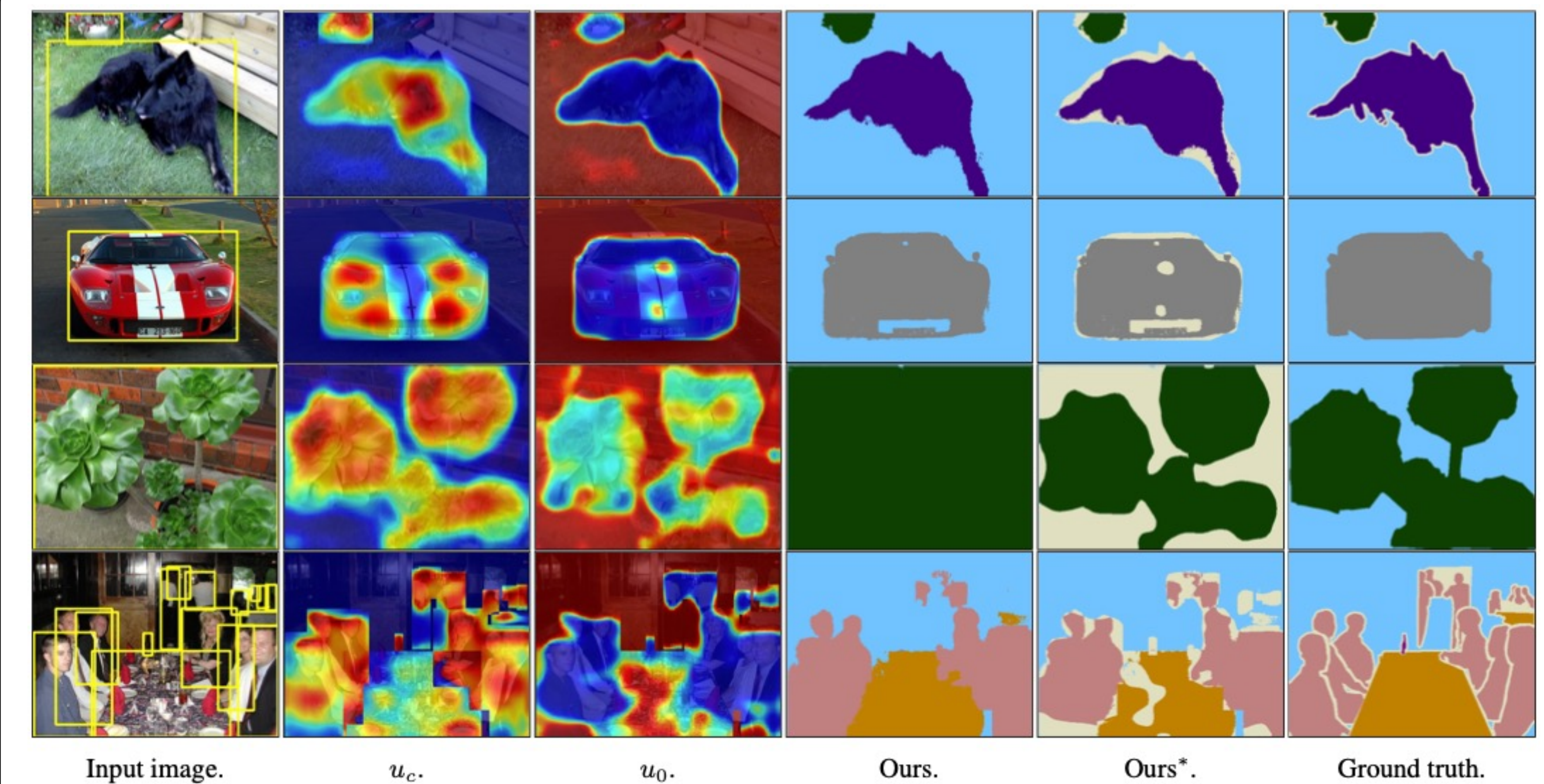
$$\mathcal{L}_{wce} = -\frac{1}{\sum_c \sum_{\mathbf{p} \in \sim S_c} \sigma(\mathbf{p})} \sum_c \sum_{\mathbf{p} \in \sim S_c} \sigma(\mathbf{p}) \log H_c(\mathbf{p})$$

$$D_c(\mathbf{p}) = 1 + \left( \frac{\phi(\mathbf{p})}{\|\phi(\mathbf{p})\|} \cdot \frac{W_c}{\|W_c\|} \right)$$

$$\sigma(\mathbf{p}) = \left( \frac{D_{c^*}(\mathbf{p})}{\max_c(D_c(\mathbf{p}))} \right)^\tau$$

$\sigma$ : a confidence map where values approach to one when the label  $c^* = Y_{crf}(\mathbf{p})$  is confident  
 $\sim S_c$ : a set of locations where  $Y_{crf}$  and  $Y_{ret}$  give different labels

## Experiments



Method	val	test
<i>Supervision: Image-level labels (10K) with Saliency (3K)</i>		
SeeNet (NIPS 2018)	61.1	60.7
FickleNet (CVPR 2019)	61.2	61.9
OAA (ICCV 2019)	63.1	62.8
ICD (CVPR 2020)	64.0	63.9
<i>Supervision: Boxes (10K)</i>		
BoxSup (ICCV 2015)	62.0	64.6
WSSL (ICCV 2015)	60.6	62.2
SDI (CVPR 2017)	65.7	<u>67.5</u>
BCM (CVPR 2019)	66.8	-
Ours		
w/ $Y_{crf}$	<u>67.8</u>	-
w/ $Y_{ret}$	66.1	-
w/ NAL	<b>68.1</b>	<b>69.4</b>
<i>Supervision: Boxes (9K) with Masks (1K)</i>		
BoxSup (ICCV 2015)	63.5	66.2
WSSL (ICCV 2015)	65.1	66.6
SDI (CVPR 2017)	65.8	<u>66.9</u>
BCM (CVPR 2019)	<u>67.5</u>	-
Ours w/ NAL	<b>70.5</b>	<b>71.5</b>

mIoUs of DeepLab-V1 on PASCAL VOC 2012