

# Learning Multi-modal Class-specific Tokens for Weakly Supervised Dense Object Localization

## -Supplementary Material-

### A. Implementation details

#### Classification training and dense localization inference.

For the training of the proposed multi-modal token transformer, publicly available pre-trained weights<sup>1</sup> were used to initialize the ViT-base backbone. The decay parameter of the Global Weighted Ranking Pooling (GWRP) was set to 0.996, as suggested in [6]. We selected the model with the best classification performance on the validation set for the dense localization inference. The dense localization inference process is illustrated in Figure 1. For the generation of the transformer attention-driven localization maps, we fused the class-to-patch attention maps from the last six transformer encoding layers. The patch-level pairwise affinity maps were generated by fusing the patch-to-patch attention maps from all twelve transformer encoding layers. Following the prior works [5, 9, 13, 14], we reported the best evaluation results of the dense object localization maps when applying multiple background thresholds.

**Segmentation training and inference.** As in previous works [4, 7, 11, 12, 14], we used ResNet38-based Deeplab-V1 as the semantic segmentation network. We first trained the ResNet38-based classification network for 15 epochs using the image-level labels of the target segmentation dataset, following the same settings as [2]. We then used the classification weights to initialize the segmentation network. During training, we first randomly scaled the images within the range of (0.7, 1.3) and randomly applied horizontal flipping on the scaled images. Moreover, we randomly cropped the processed images to be of size  $321 \times 321$ . The initial learning rate was set as  $7 \times 10^{-4}$ . A polynomial learning rate decay with a power of 0.9 was used. With the stochastic gradient descent (SGD) optimizer, we trained the segmentation model for 30 epochs with a batch size of 4. For inference, we used multi-scale inputs of scales 0.75, 1.0 and 1.5. Results were max-pooled and post-processed by the CRF with the default hyper-parameters suggested in [3]. The same settings were used on both PASCAL VOC 2012 and MS COCO 2014.

### B. Additional quantitative results

Following prior works [4, 5, 8–10], we used IRN [1] to post-process our generated class-specific dense localization maps (seeds), generating the pseudo masks (masks) for WSSS. Table 1 shows the evaluation results of the seeds and

<sup>1</sup>[https://dl.fbaipublicfiles.com/dino/dino\\_vitbasel6\\_pretrain/dino\\_vitbasel6\\_pretrain.pth](https://dl.fbaipublicfiles.com/dino/dino_vitbasel6_pretrain/dino_vitbasel6_pretrain.pth)

Table 1. Evaluation of the generated multi-label dense localization maps (seed) and their post-processed masks using IRN [1] for WSSS in terms of mIoU (%) on the PASCAL VOC 2012 *train* set.

Method	Cls. backbone	Seed	Mask
IRN (CVPR19) [1]	ResNet50	48.8	66.3
CONTA (NeurIPS20) [15]	ResNet50	48.8	67.9
RIB (NeurIPS21) [8]	ResNet50	56.5	70.6
AdvCAM (CVPR21) [9]	ResNet38	55.6	69.9
CDA (ICCV21) [11]	ResNet38	50.8	67.7
Du <i>et al.</i> (CVPR22) [5]	ResNet38	61.5	70.1
W-OoD (CVPR22) [10]	ResNet50	59.1	72.1
Ours	ViT-base	<b>66.3</b>	<b>73.7</b>

masks generated by the state-of-the-art methods. The proposed method is seen to achieve the best mIoUs in terms of both seeds and pseudo masks. Table 2 and Table 3 present the detailed segmentation results of per class IoUs on PASCAL VOC 2012 and MS COCO 2014, respectively.

### C. Additional qualitative results

**Qualitative dense object localization results.** Figure 3 and Figure 4 present the object localization heatmaps (where the colors from red to blue indicate the activation scores from high to low) generated by the proposed method on the multi-label PASCAL VOC 2012 and MS COCO 2014 *train* sets, respectively. Compared to the results by the state-of-the-art method, i.e., MCTformer [14], the proposed method is shown to produce more complete object localization maps, even for the challenging cases, such as the multiple TV monitors in a complex background in the last row of Figure 3 and the person in a far corner of the scene in the last row of Figure 4. Figure 5 presents the object localization heatmaps generated by the proposed method on the single-label OpenImages *test* set. Table 5 shows that the proposed method produces accurate dense object localization maps with clear boundaries, where the object regions are generally activated with high scores. These results demonstrate the effectiveness of the proposed method in performing dense object localization on both multi-label and single-label images.

**Qualitative segmentation results.** Figure 6 and Figure 7 present additional qualitative segmentation results on the PASCAL VOC 2012 and MS COCO 2014 *val* sets, respectively. Using the pseudo labels generated by the proposed method, the trained segmentation models are seen to perform well on both PASCAL VOC and MS COCO in vari-

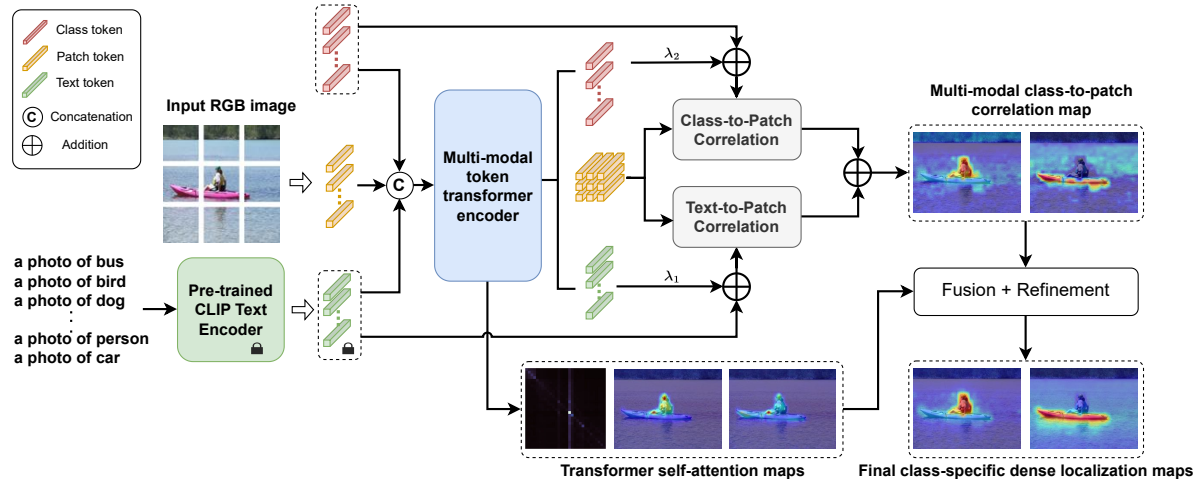


Figure 1. Class-specific dense localization inference of the proposed method.

Table 2. Segmentation performance comparison with the state-of-the-art WSSS methods using only image-level labels in terms of per-class segmentation IoUs (%) on PASCAL VOC 2012. \* denotes without post-processing (*i.e.*, multi-scale testing and CRF).

	bkg	plane	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbk	person	plant	sheep	sofa	train	tv	mIoU
Results on the <i>val</i> set:																						
AdvCAM [9]	90.0	79.8	34.1	82.6	63.3	70.5	<b>89.4</b>	76.0	87.3	31.4	81.3	33.1	82.5	80.8	74.0	72.9	50.3	82.3	42.2	74.1	52.9	68.1
Zhang <i>et al.</i> [16]	89.9	75.1	32.9	87.8	60.9	69.5	87.7	79.5	89.0	28.0	80.9	34.8	83.4	79.7	74.7	66.9	<b>56.5</b>	82.7	44.9	73.1	45.7	67.8
MCTformer [14]	91.9	78.3	39.5	<b>89.9</b>	55.9	76.7	81.8	79.0	90.7	32.6	87.1	57.2	87.0	84.6	77.4	79.2	55.1	<b>89.2</b>	47.2	70.4	58.8	71.9
W-OoD [10]	91.0	80.1	34.1	88.1	<b>64.8</b>	68.3	87.4	<b>84.4</b>	<b>89.8</b>	30.1	<b>87.8</b>	34.7	<b>87.5</b>	<b>85.9</b>	<b>79.8</b>	75.0	56.4	84.5	<b>47.8</b>	<b>80.4</b>	46.4	70.7
Ours*	91.0	78.5	38.2	85.3	58.8	73.3	84.7	80.5	85.4	26.4	79.8	56.0	79.3	77.6	73.7	78.5	47.8	82.6	44.0	68.4	56.8	68.9
Ours	<b>92.4</b>	<b>84.7</b>	<b>42.2</b>	85.5	64.1	<b>77.4</b>	86.6	82.2	88.7	<b>32.7</b>	83.8	<b>59.0</b>	82.4	80.9	76.1	<b>81.4</b>	48.0	88.2	46.4	70.2	<b>62.5</b>	<b>72.2</b>
Results on the <i>test</i> set:																						
AdvCAM [9]	90.1	81.2	33.6	80.4	52.4	66.6	87.1	80.5	87.2	28.9	80.1	38.5	84.0	83.0	79.5	71.9	47.5	80.8	59.1	65.4	49.7	68.0
Zhang <i>et al.</i> [16]	90.4	79.8	32.9	85.8	52.9	66.4	<b>87.2</b>	<b>81.4</b>	87.6	28.2	79.7	50.2	82.9	80.4	78.9	70.6	51.2	83.4	55.4	68.5	44.6	68.5
MCTformer [14]	92.3	84.4	37.2	82.8	60.0	72.8	78.0	79.0	89.4	31.7	<b>84.5</b>	<b>59.1</b>	<b>85.3</b>	<b>83.8</b>	79.2	81.0	53.9	85.3	60.5	65.7	57.7	71.6
W-OoD [10]	90.9	83.1	35.6	<b>89.0</b>	<b>61.5</b>	63.0	86.2	80.8	<b>89.9</b>	29.6	79.6	40.1	82.1	81.0	<b>82.6</b>	74.0	<b>60.1</b>	85.3	58.0	<b>71.9</b>	47.0	70.1
Ours*	91.2	82.8	36.6	81.5	54.6	70.2	81.6	79.2	83.9	31.2	78.0	57.8	80.8	81.6	77.6	78.4	50.4	80.9	57.3	58.5	58.7	69.2
Ours	<b>92.6</b>	<b>87.9</b>	<b>40.4</b>	86.2	56.7	<b>75.3</b>	82.5	81.1	85.9	<b>33.8</b>	82.4	57.8	84.6	83.1	80.0	<b>81.1</b>	52.8	<b>86.2</b>	<b>61.5</b>	61.1	<b>63.8</b>	<b>72.2</b>

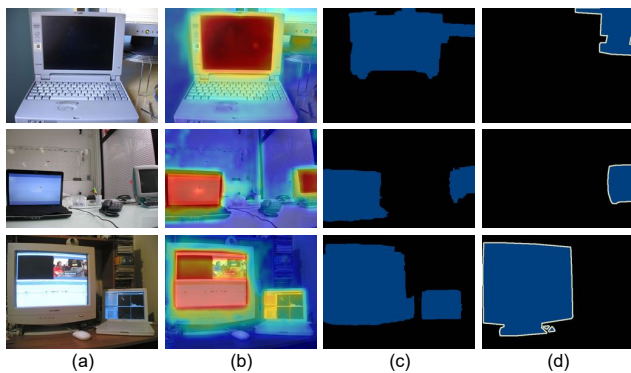


Figure 2. Failure cases. (a) Input. (b) The generated dense localization maps for the class “TV monitor”. (c) The corresponding pseudo labels. (d) Ground-truth semantic segmentation labels.

ous challenging indoor and outdoor scenarios, such as small objects and complex backgrounds. This demonstrates the effectiveness and the good generalization ability of the proposed method in WSSS.

**Limitation discussion.** Our method has a limitation in accurately distinguishing between the target objects and some background objects that have very similar visual features. For example, as shown in Figure 2, the generated dense localization maps for the class “TV monitor” mistakenly include the regions of the laptop given that they are similar in shape and both have screens.

Table 3. Segmentation performance comparison with the state-of-the-art WSSS methods in terms of per class IoU(%) on the MS COCO 2014 *val* set. \* denotes without post-processing.

Class	MCTformer [14]	Our*	Ours	Class	MCTformer [14]	Our*	Ours
background	82.4	84.2	<b>85.3</b>	wine class	27.0	<b>35.3</b>	33.8
person	62.6	71.2	<b>72.9</b>	cup	29.0	33.6	<b>35.8</b>
bicycle	47.4	<b>50.8</b>	49.8	fork	13.9	18.7	<b>20.0</b>
car	<b>47.2</b>	47.0	43.8	knife	12.0	10.8	<b>12.6</b>
motorcycle	63.7	65.1	<b>66.2</b>	spoon	6.6	<b>6.7</b>	<b>6.7</b>
airplane	64.7	64.4	<b>69.2</b>	bowl	22.4	<b>24.0</b>	23.7
bus	64.5	67.0	<b>69.1</b>	banana	63.2	61.4	<b>64.4</b>
train	<b>64.5</b>	62.3	63.7	apple	44.4	45.4	<b>50.8</b>
truck	<b>44.8</b>	40.3	43.4	sandwich	39.7	44.3	<b>47.0</b>
boat	<b>42.3</b>	41.1	<b>42.3</b>	orange	63.0	60.6	<b>64.6</b>
traffic light	<b>49.9</b>	48.8	49.3	broccoli	<b>51.2</b>	50.5	50.6
fire hydrant	73.2	72.1	<b>74.9</b>	carrot	<b>40.0</b>	34.1	38.6
stop sign	76.6	72.6	<b>77.3</b>	hot dog	53.0	51.0	<b>54.0</b>
parking meter	64.4	62.9	<b>67.0</b>	pizza	62.2	63.3	<b>64.1</b>
bench	32.8	<b>35.2</b>	34.1	donut	55.7	55.4	<b>59.7</b>
bird	62.6	62.1	<b>63.1</b>	cake	47.9	47.1	<b>50.6</b>
cat	<b>78.2</b>	74.3	76.2	chair	22.8	<b>24.5</b>	<b>24.5</b>
dog	68.2	67.3	<b>70.6</b>	couch	35.0	37.7	<b>40.0</b>
horse	65.8	64.7	<b>67.1</b>	potted plant	13.5	<b>16.8</b>	13.0
sheep	70.1	68.0	<b>70.8</b>	bed	48.6	51.5	<b>53.7</b>
cow	68.3	67.0	<b>71.2</b>	dining table	12.9	<b>24.7</b>	19.2
elephant	81.6	79.9	<b>82.2</b>	toilet	63.1	64.6	<b>66.6</b>
bear	<b>80.1</b>	76.2	79.6	tv	47.9	47.6	<b>50.8</b>
zebra	<b>83.0</b>	82.0	82.8	laptop	49.5	52.8	<b>55.4</b>
giraffe	<b>76.9</b>	76.4	76.7	mouse	13.4	11.8	<b>14.4</b>
backpack	14.6	16.7	<b>17.5</b>	remote	41.9	40.7	<b>47.1</b>
umbrella	61.7	60.6	<b>66.9</b>	keyboard	49.8	53.3	<b>57.2</b>
handbag	4.5	<b>8.1</b>	5.8	cellphone	54.1	50.8	<b>54.9</b>
tie	25.2	29.3	<b>31.4</b>	microwave	38.0	42.7	<b>46.1</b>
suitcase	46.8	47.4	<b>51.4</b>	oven	29.9	<b>35.4</b>	35.3
frisbee	43.8	46.1	<b>54.1</b>	toaster	0.0	<b>4.3</b>	2.0
skis	12.8	11.6	<b>13.0</b>	sink	28.0	29.3	<b>36.1</b>
snowboard	<b>31.4</b>	26.6	30.3	refrigerator	40.1	50.2	<b>52.7</b>
sports ball	9.2	32.3	<b>36.1</b>	book	32.2	31.2	<b>34.8</b>
kite	26.3	36.2	<b>47.5</b>	clock	43.2	46.3	<b>51.5</b>
baseball bat	0.9	<b>7.2</b>	7.0	vase	22.6	24.4	<b>25.8</b>
baseball globe	0.7	8.5	<b>10.4</b>	scissors	<b>32.9</b>	29.9	30.7
skateboard	7.8	11.7	<b>15.2</b>	teddy bear	<b>61.9</b>	60.0	61.4
surfboard	46.5	45.2	<b>51.5</b>	hair drier	0.0	<b>1.8</b>	1.3
tennis racket	1.4	21.2	<b>26.4</b>	toothbrush	12.2	16.9	<b>19.0</b>
bottle	31.1	36.1	<b>37.1</b>	<b>mIoU</b>	42.0	43.7	<b>45.9</b>

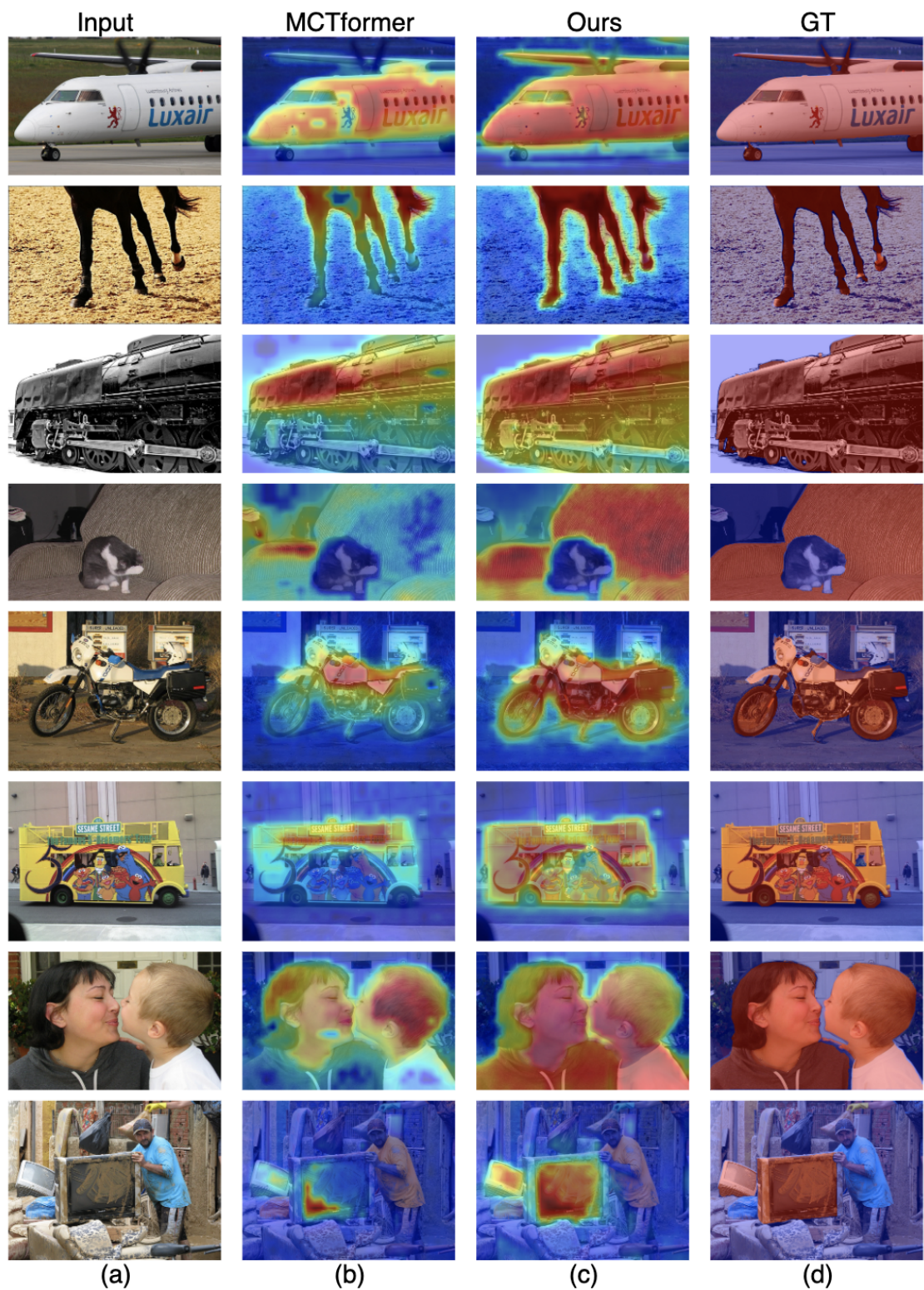


Figure 3. Visualization of the object localization heatmaps on the PASCAL VOC 2012 *train* set. (a) Input. (b) Results of MCTformer [14]. (c) Our results. (d) Ground-truth. (Only the localization heatmaps of the dominant class in each image are presented.)

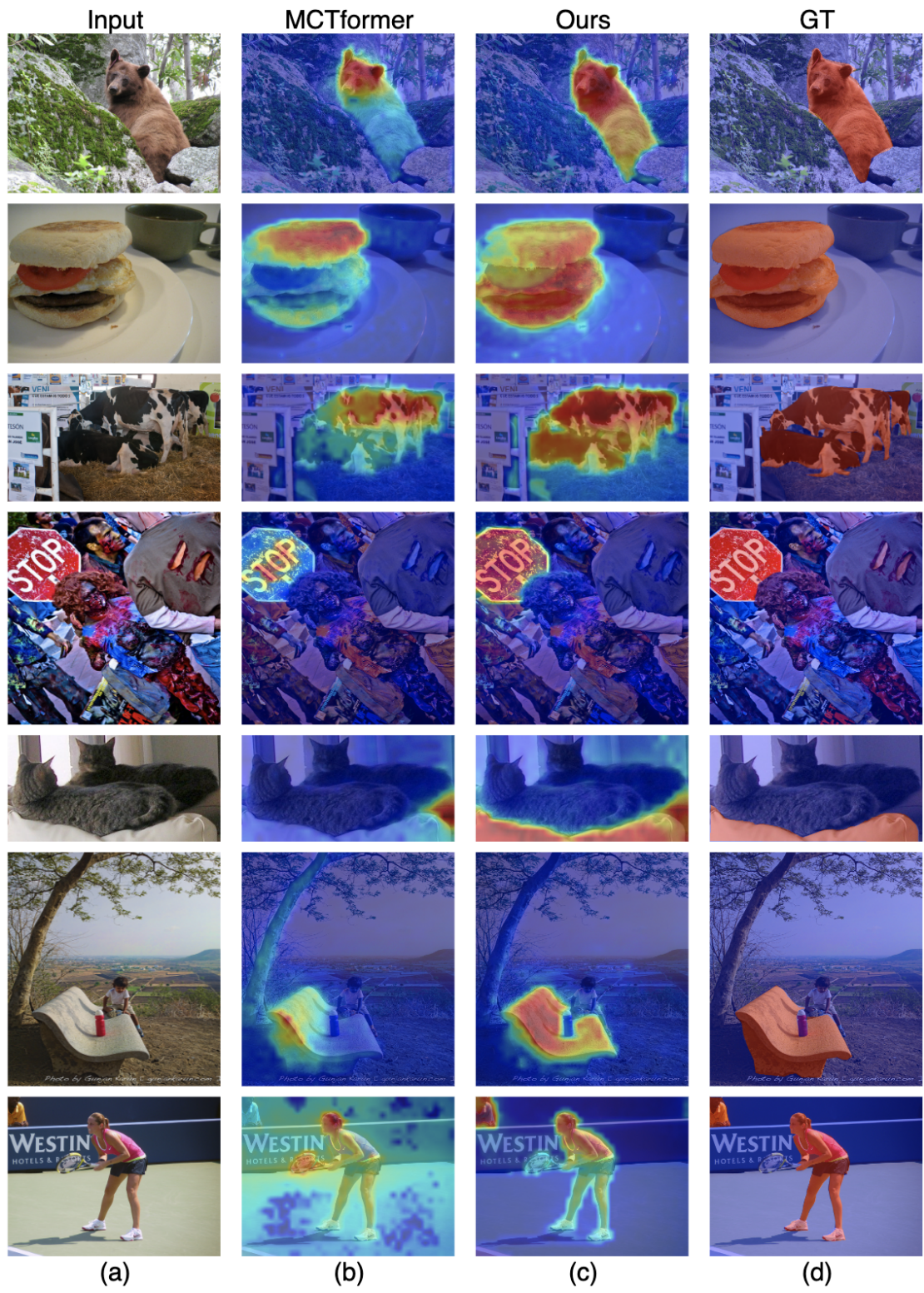


Figure 4. Visualization of the object localization heatmaps on the MS COCO 2014 *train* set. (a) Input. (b) Results of MCTformer [14]. (c) Our results. (d) Ground-truth. (Only the localization heatmaps of the dominant class in each image are presented.)

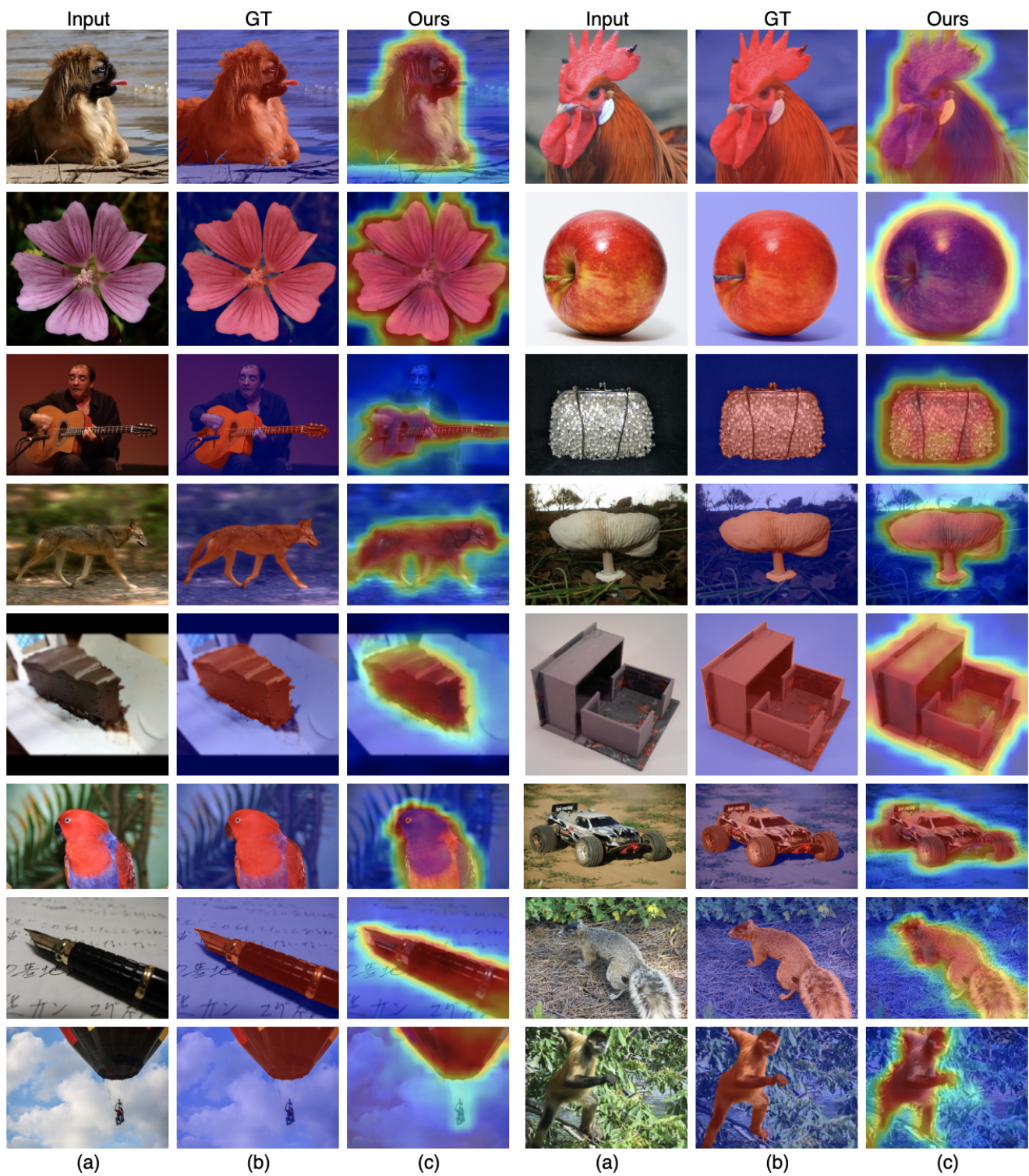


Figure 5. Visualization of the object localization heatmaps on the single-label OpenImages *test* set. (a) Input. (b) Ground-truth. (c) Our results.

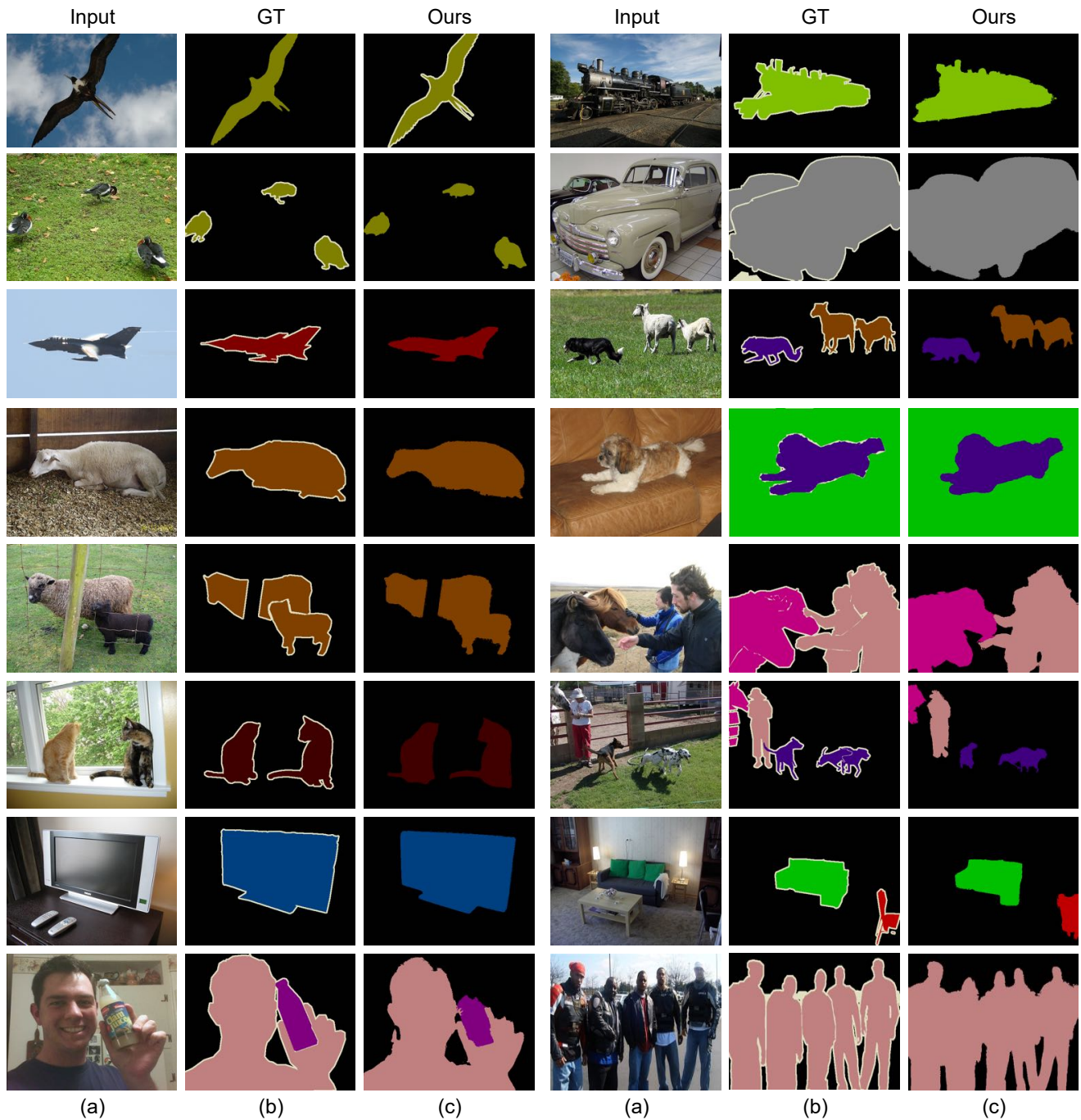


Figure 6. Visualization of the semantic segmentation results on the PASCAL VOC 2012 *val* set. (a) Input. (b) Ground-truth. (c) Our results.

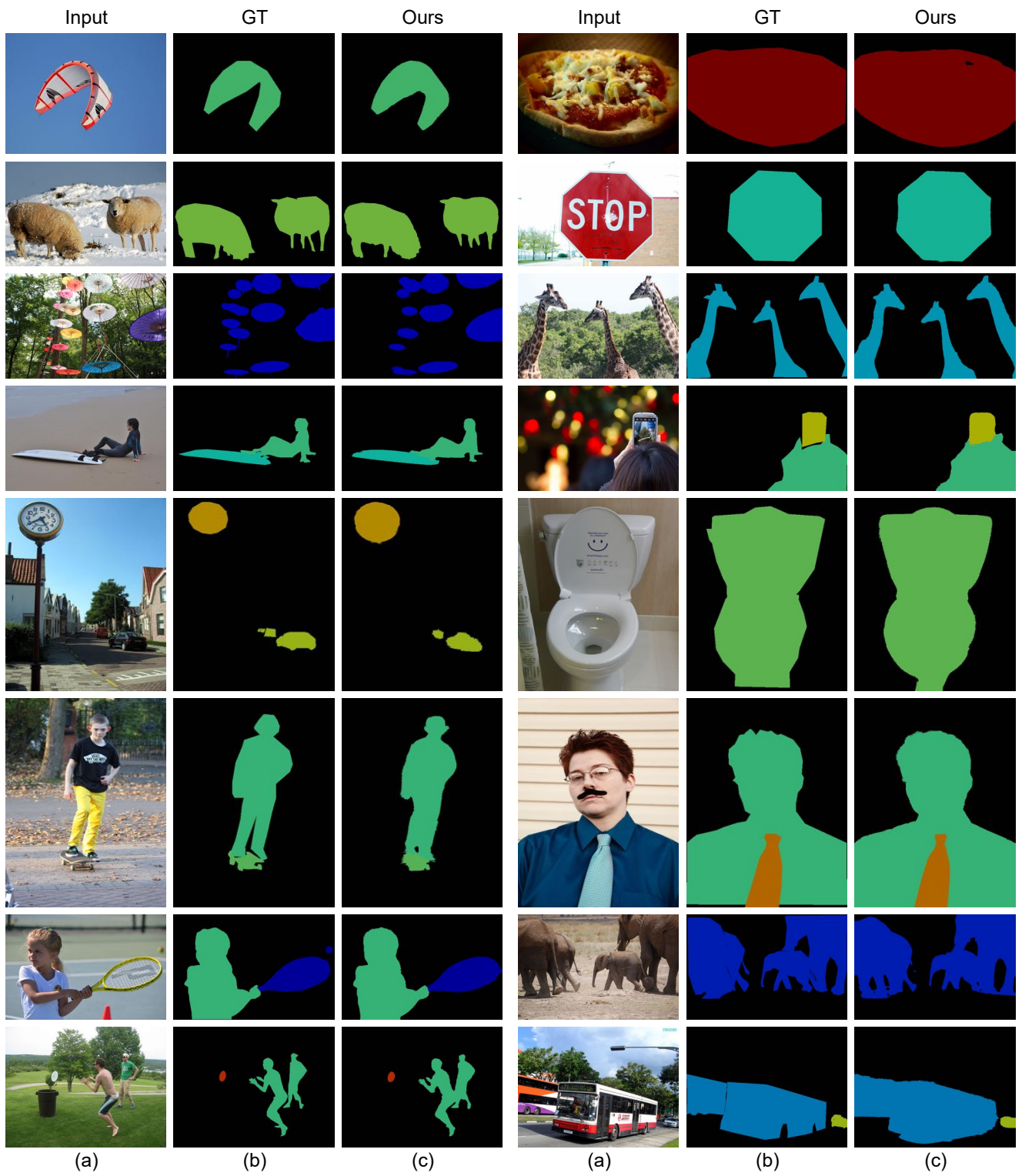


Figure 7. Visualization of the semantic segmentation results on the MS COCO 2014 *val* set. (a) Input. (b) Ground-truth. (c) Our results.



## References

- [1] Jiwoon Ahn, Sunghyun Cho, and Suha Kwak. Weakly supervised learning of instance segmentation with inter-pixel relations. In *CVPR*, 2019. 1
- [2] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *CVPR*, 2018. 1
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 1
- [4] Qi Chen, Lingxiao Yang, Jian-Huang Lai, and Xiaohua Xie. Self-supervised image-specific prototype exploration for weakly supervised semantic segmentation. In *CVPR*, 2022. 1
- [5] Ye Du, Zehua Fu, Qingjie Liu, and Yunhong Wang. Weakly supervised semantic segmentation by pixel-to-prototype contrast. In *CVPR*, 2022. 1
- [6] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *ECCV*, 2016. 1
- [7] Hyeokjun Kweon, Sung-Hoon Yoon, Hyeonseong Kim, Daehee Park, and Kuk-Jin Yoon. Unlocking the potential of ordinary classifier: Class-specific adversarial erasing framework for weakly supervised semantic segmentation. In *ICCV*, 2021. 1
- [8] Jungbeom Lee, Jooyoung Choi, Jisoo Mok, and Sungroh Yoon. Reducing information bottleneck for weakly supervised semantic segmentation. *NeurIPS*, 2021. 1
- [9] Jungbeom Lee, Eunji Kim, and Sungroh Yoon. Anti-adversarially manipulated attributions for weakly and semi-supervised semantic segmentation. In *CVPR*, 2021. 1, 2
- [10] Jungbeom Lee, Seong Joon Oh, Sangdoon Yun, Junsuk Choe, Eunji Kim, and Sungroh Yoon. Weakly supervised semantic segmentation using out-of-distribution data. In *CVPR*, 2022. 1, 2
- [11] Yukun Su, Ruizhou Sun, Guosheng Lin, and Qingyao Wu. Context decoupling augmentation for weakly supervised semantic segmentation. In *ICCV*, 2021. 1
- [12] Kunyang Sun, Haoqing Shi, Zhengming Zhang, and Yongming Huang. Ecs-net: Improving weakly supervised semantic segmentation by using connections between class activation maps. In *ICCV*, 2021. 1
- [13] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *CVPR*, 2020. 1
- [14] Lian Xu, Wanli Ouyang, Mohammed Bennamoun, Farid Boussaid, and Dan Xu. Multi-class token transformer for weakly supervised semantic segmentation. In *CVPR*, 2022. 1, 2, 3, 4, 5
- [15] Dong Zhang, Hanwang Zhang, Jinhui Tang, Xiansheng Hua, and Qianru Sun. Causal intervention for weakly-supervised semantic segmentation. In *NeurIPS*, 2020. 1
- [16] Fei Zhang, Chaochen Gu, Chenyue Zhang, and Yuchao Dai. Complementary patch for weakly supervised semantic segmentation. In *ICCV*, 2021. 2