

# Unseen Object Segmentation in Videos via Transferable Representations

Yi-Wen Chen<sup>1,2</sup>, Yi-Hsuan Tsai<sup>3</sup>, Chu-Ya Yang<sup>1</sup>,  
Yen-Yu Lin<sup>1</sup>, Ming-Hsuan Yang<sup>4,5</sup>

<sup>1</sup>Academia Sinica, <sup>2</sup>National Taiwan University, <sup>3</sup>NEC Laboratories America,  
<sup>4</sup>University of California, Merced, <sup>5</sup>Google Cloud

## 1 Analysis of Transferring Visual Information

We analyze the proposed method for transferring visual information by investigating the weights of the transferable layer. Table 1 presents the learned weights of the transferable layer on the DAVIS dataset for unseen object categories. For each target video, the source categories with higher weights are similar to the target video in appearance, which gives reasonable transform of visual information.

**Table 1.** Learned weights of the transferable layer on the DAVIS dataset for transferring knowledge from seen/source categories (rows) to unseen/target object categories (columns). For each unseen category, the largest weight over all seen categories is marked in bold.

Sequence	bear	bswan	camel	eleph	goat	malw	rhino
aero	0.286	0.419	0.381	0.412	0.279	0.430	0.325
bike	0.317	0.372	0.393	0.423	0.358	0.309	0.432
bird	0.624	<b>0.891</b>	0.538	0.572	0.614	<b>0.780</b>	0.595
boat	0.392	0.419	0.358	0.460	0.323	0.474	0.428
bottle	0.401	0.336	0.307	0.410	0.349	0.387	0.368
bus	0.392	0.262	0.266	0.440	0.306	0.200	0.327
car	0.488	0.317	0.469	0.559	0.379	0.292	0.508
cat	<b>0.756</b>	0.436	0.417	0.574	0.609	0.398	0.492
chair	0.507	0.314	0.406	0.528	0.466	0.362	0.450
cow	0.701	0.409	0.715	0.748	0.618	0.346	<b>0.846</b>
table	0.341	0.310	0.186	0.301	0.291	0.504	0.257
dog	0.700	0.476	0.534	0.603	<b>0.788</b>	0.417	0.576
horse	0.547	0.330	<b>0.898</b>	<b>0.770</b>	0.692	0.260	0.776
mbike	0.301	0.287	0.346	0.408	0.371	0.287	0.355
person	0.504	0.429	0.731	0.639	0.554	0.366	0.629
plant	0.463	0.418	0.364	0.437	0.428	0.451	0.474
sheep	0.721	0.525	0.491	0.662	0.616	0.348	0.605
sofa	0.366	0.309	0.366	0.447	0.404	0.291	0.412
train	0.298	0.260	0.343	0.488	0.320	0.204	0.419
tv	0.369	0.252	0.277	0.425	0.271	0.248	0.303

## 2 Runtime Performance

The runtime performance is shown in Table 2. All the timings are measured on a machine with 2.5 GHz Intel Xeon CPU. We compute the optical flow [7] and utilize the minimum barrier distance [13] to generate motion prior using MATLAB implementations. In the proposed formulation, feature extraction, segment mining, and CNN model training are implemented using Python and Caffe library on a GPU of NVIDIA GTX 1080 Ti with 11 GB memory. The CNN model is fine-tuned for 2000 iterations. Note that, we report the timings for each component during iterative optimization averaged on all the frames.

**Table 2.** Runtime performance on the DAVIS dataset.

Stage	Time (second)
Motion prior computing (per pair of frames)	0.01
Feature extraction (per frame)	1.72
Segment mining (per frame)	0.01
CNN model training (per frame)	7.31

## 3 Per-video Results on the DAVIS 2016 Dataset

In Table 3, we present the results of each video from the DAVIS 2016 dataset under weakly-supervised and unsupervised settings. We show that the proposed algorithm achieves better performance than the state-of-the-art methods in most videos.

**Table 3.** Per-video results on the DAVIS 2016 dataset.

Methods	Weak Supervision			No Supervision				
	SPFTN [12]	FCN [8]	Ours	MSG [10]	FST [9]	NLC [1]	FSEG [4]	Ours
bear	74.8	80.3	<b>89.8</b>	85.1	89.8	90.7	91.5	<b>91.8</b>
bswan	<b>87.6</b>	75.6	76.7	52.6	73.2	87.5	89.5	<b>90.3</b>
bumps	29.7	29.9	<b>36.2</b>	35.3	24.1	<b>63.5</b>	38.8	42.1
trees	35.0	29.2	<b>40.5</b>	18.8	18.0	21.2	34.7	<b>38.9</b>
boat	35.9	63.4	<b>67.0</b>	14.4	36.1	0.7	63.8	<b>63.8</b>
bdan	37.1	14.6	<b>46.0</b>	23.6	46.7	<b>67.3</b>	14.2	13.1
bdanF	70.0	51.4	<b>80.0</b>	15.7	61.6	<b>80.4</b>	54.9	62.7
bus	<b>81.5</b>	61.1	81.2	<b>88.5</b>	82.5	62.9	80.4	80.5
camel	<b>76.2</b>	70.9	72.0	75.6	56.2	76.8	76.4	<b>77.5</b>
carR	76.8	71.0	<b>88.8</b>	63.0	<b>80.8</b>	50.9	74.8	79.6
carS	78.1	87.1	<b>92.5</b>	88.0	69.8	64.5	88.4	<b>93.3</b>

carT	75.4	86.7	<b>90.4</b>	62.1	85.1	83.3	90.7	<b>92.5</b>
cows	77.0	85.7	<b>88.1</b>	79.9	79.1	88.3	88.0	<b>88.3</b>
jump	34.2	33.6	<b>63.8</b>	6.5	59.8	<b>71.8</b>	10.3	11.2
twirl	46.1	27.8	<b>65.5</b>	36.6	45.3	34.7	<b>46.2</b>	41.0
dog	85.6	71.2	<b>89.1</b>	33.1	70.8	80.9	90.4	<b>91.6</b>
dogA	7.1	39.3	<b>72.9</b>	11.0	28.0	65.2	<b>68.9</b>	65.1
drtC	55.9	58.9	<b>67.1</b>	<b>75.8</b>	66.7	32.4	46.1	65.1
drtS	62.3	69.9	<b>79.4</b>	57.5	<b>68.3</b>	47.3	67.2	66.4
drtT	67.8	76.4	<b>80.6</b>	63.8	53.3	15.4	85.1	<b>89.7</b>
eleph	<b>75.6</b>	70.4	73.8	68.9	82.4	51.8	<b>86.2</b>	85.7
flamg	<b>38.1</b>	33.5	34.5	79.4	<b>81.7</b>	53.9	44.5	47.8
goat	72.8	83.1	<b>83.3</b>	73.5	55.4	1.0	84.1	<b>84.8</b>
hike	<b>89.3</b>	84.1	79.0	60.3	88.9	<b>91.8</b>	82.5	83.4
hockey	60.2	72.7	<b>73.1</b>	71.3	46.7	<b>81.0</b>	66.0	70.7
hjH	35.1	<b>77.6</b>	67.0	73.4	57.8	<b>83.4</b>	71.1	72.1
hjL	41.1	<b>79.5</b>	73.6	68.2	52.6	65.1	70.2	<b>76.5</b>
ksurf	<b>58.3</b>	55.8	46.5	41.9	27.2	45.3	47.7	<b>49.0</b>
kwalk	<b>73.3</b>	52.1	48.9	59.7	64.9	<b>81.3</b>	52.7	51.3
libby	50.8	49.5	<b>59.4</b>	5.0	50.7	63.5	67.7	<b>68.1</b>
lucia	83.3	<b>84.2</b>	78.9	41.7	64.4	<b>87.6</b>	79.9	81.0
malf	<b>70.8</b>	47.5	45.8	3.3	60.1	61.7	74.6	<b>75.2</b>
malw	<b>65.8</b>	40.9	41.6	4.5	8.7	76.1	83.3	<b>84.9</b>
motob	75.0	<b>77.7</b>	71.6	46.6	61.7	61.4	83.8	<b>85.2</b>
motoj	60.8	61.5	<b>65.5</b>	61.8	60.2	25.1	<b>80.4</b>	77.2
mbike	47.6	<b>78.5</b>	58.4	<b>73.8</b>	55.9	71.4	28.7	38.6
parag	<b>72.6</b>	30.9	28.1	<b>93.3</b>	72.5	88.0	17.7	5.5
paral	<b>62.8</b>	57.0	58.1	51.2	50.6	<b>62.8</b>	58.9	59.4
park	67.7	<b>84.0</b>	78.2	29.5	45.8	<b>90.1</b>	79.4	79.5
rhino	55.2	57.7	<b>71.0</b>	<b>90.2</b>	77.6	68.2	77.6	86.0
rolb	12.5	64.2	<b>73.2</b>	80.1	31.8	<b>81.4</b>	63.3	72.7
scbla	58.8	45.0	<b>72.1</b>	<b>57.9</b>	52.2	16.2	36.1	36.4
scgra	67.0	<b>73.7</b>	72.9	34.5	32.5	58.7	73.2	<b>75.7</b>
sobox	<b>57.8</b>	47.5	51.9	<b>67.2</b>	41.0	63.4	49.7	47.4
socB	49.0	<b>49.5</b>	46.3	37.0	<b>84.3</b>	82.9	29.3	28.3
strol	<b>65.4</b>	58.7	58.7	<b>67.8</b>	58.0	84.9	63.9	62.8
surf	<b>87.0</b>	78.4	79.1	77.0	47.5	77.5	88.8	<b>91.2</b>
swing	75.5	75.5	<b>76.4</b>	62.2	43.1	<b>85.1</b>	73.8	74.0
tennis	62.5	<b>78.2</b>	73.0	59.0	38.8	<b>87.1</b>	76.9	78.4
train	73.6	46.9	<b>77.3</b>	<b>88.7</b>	83.1	72.9	42.5	51.1
Avg.	61.2	61.6	<b>67.7</b>	54.3	57.5	64.1	64.7	<b>66.5</b>

## 4 Results on the SegTrack v2 Dataset

In Table 4, we provide experiments on the SegTrack v2 dataset [6] that contains many unseen objects. We use the ResNet-101 architecture and the training data from PASCAL VOC, which is the same setting as the appearance stream in FSEG [4]. We show that the proposed method performs better than FSEG [4], other unsupervised algorithms [9,5] and HVS [2] which includes human annotations in the procedure.

**Table 4.** Results on the SegTrack v2 dataset.

Methods	FST [9]	KEY [5]	HVS [2]	FSEG [4]	Ours
Avg. IoU	53.6	57.3	50.8	56.9	58.1

## 5 Segmentation Results

We show segmentation results compared to state-of-the-art approaches on the DAVIS dataset for unseen object categories in Fig. 1-2. In the supplementary video, we present more results for each sequence with unseen categories and compare our method with baseline settings and the state-of-the-art transfer learning approach [3]. In addition, we show results using the weakly-supervised setting on the DAVIS (Fig. 3-6) and YouTube-Objects (Fig. 7-8) datasets. Some failure cases are presented in Fig. 9.

## References

1. Faktor, A., Irani, M.: Video segmentation by non-local consensus voting. In: BMVC (2014) 2
2. Grundmann, M., Kwatra, V., Han, M., Essa, I.: Efficient hierarchical graph-based video segmentation. In: CVPR (2010) 4
3. Hong, S., Oh, J., Lee, H., Han, B.: Learning transferrable knowledge for semantic segmentation with deep convolutional neural network. In: CVPR (2016) 4, 5, 6
4. Jain, S., Xiong, B., Grauman, K.: Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In: CVPR (2017) 2, 4
5. Lee, Y.J., Kim, J., Grauman, K.: Key-segments for video object segmentation. In: ICCV (2011) 4
6. Li, F., Kim, T., Humayun, A., Tsai, D., Rehg, J.M.: Video segmentation by tracking many figure-ground segments. In: ICCV (2013) 4
7. Liu, C.: Beyond pixels: Exploring new representations and applications for motion analysis. PhD thesis, MIT (2009) 2
8. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015) 2, 7, 8, 9, 10, 11, 12
9. Papazoglou, A., Ferrari, V.: Fast object segmentation in unconstrained video. In: ICCV (2013) 2, 4, 5, 6, 7, 8, 9, 10
10. P.Ochs, T.Brox: Object segmentation in video: a hierarchical variational approach for turning point trajectories into dense regions. In: ICCV (2011) 2, 5, 6, 7, 8, 9, 10



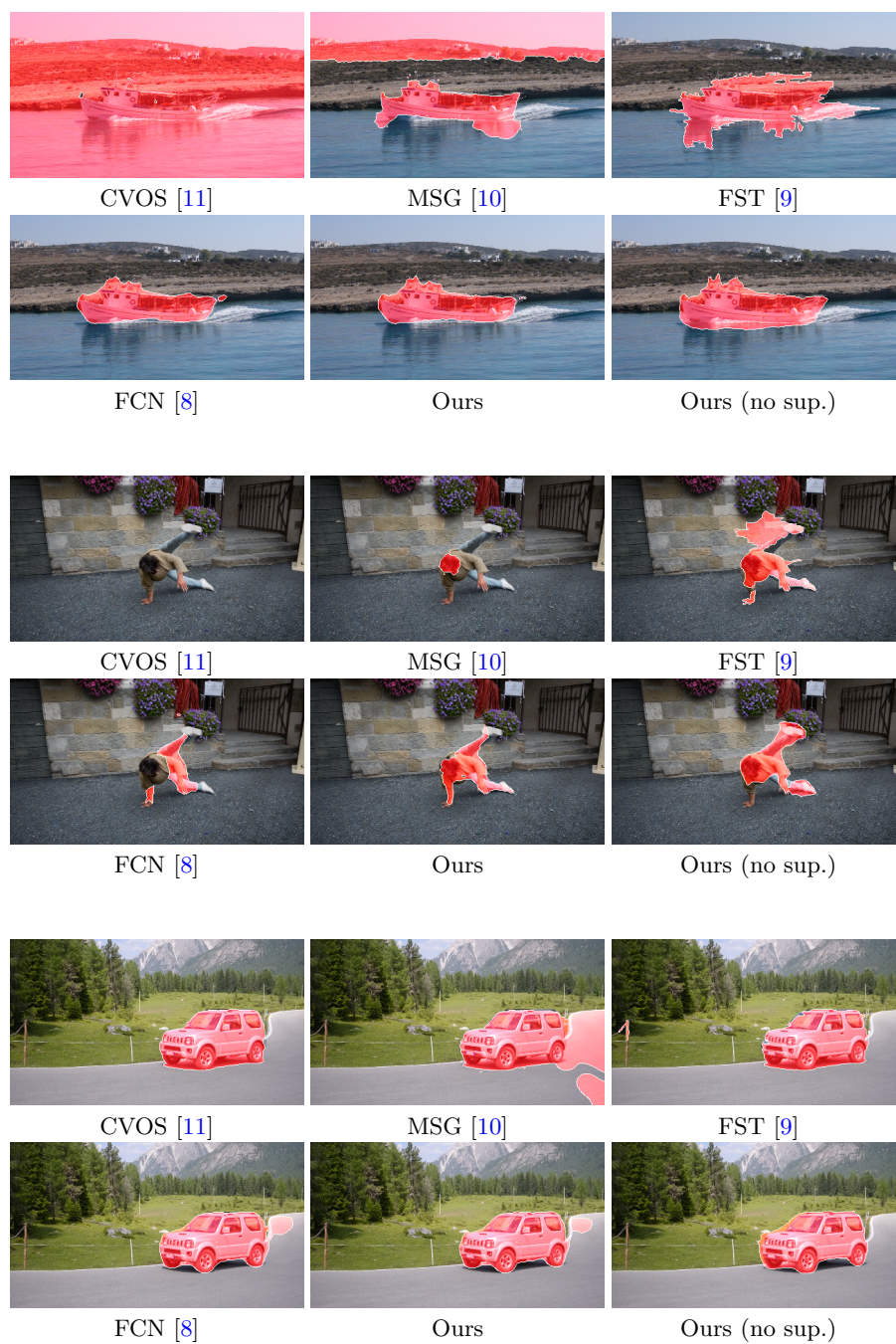
11. Taylor, B., Karasev, V., Soatto, S.: Causal video object segmentation from persistence of occlusions. In: CVPR (2015) 5, 6, 7, 8, 9, 10
12. Zhang, D., Yang, L., Meng, D., Xu, D., Han, J.: Spftn: A self-paced fine-tuning network for segmenting objects in weakly labelled videos. In: CVPR (2017) 2
13. Zhang, J., Sclaroff, S., Lin, Z., Shen, X., Price, B., Mech, R.: Minimum barrier salient object detection at 80 fps. In: ICCV (2015) 2



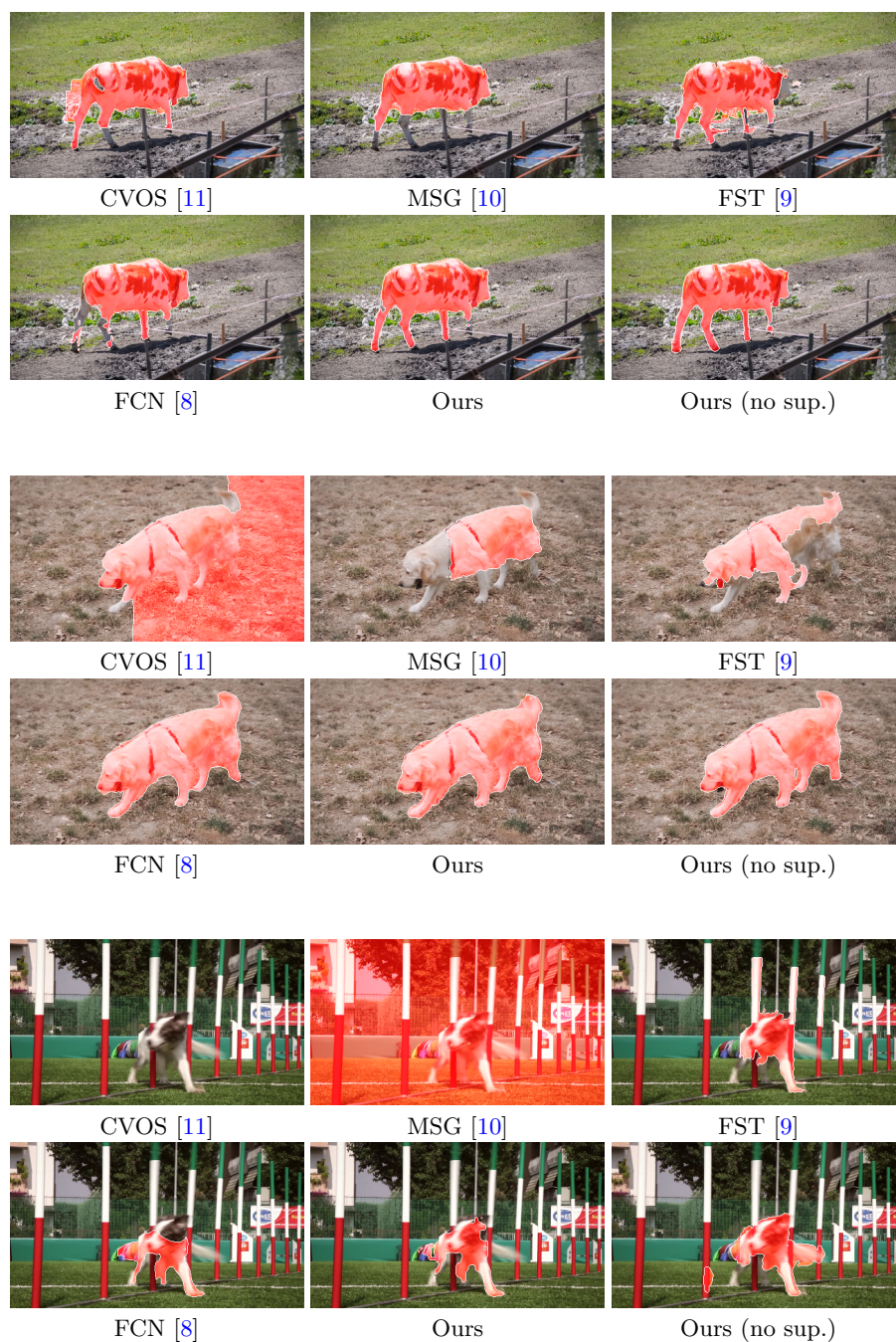
**Fig. 1.** Sample results on the DAVIS dataset for unseen object categories. Our final results contain less noisy segments and more details than other approaches and our baseline methods.



**Fig. 2.** Sample results on the DAVIS dataset for unseen object categories. Our final results contain less noisy segments and more details than other approaches and our baseline methods.



**Fig. 3.** Sample results on the DAVIS dataset with categories shared in the PASCAL VOC dataset. We show that our results with and without weak supervisions have more complete object segments with details.

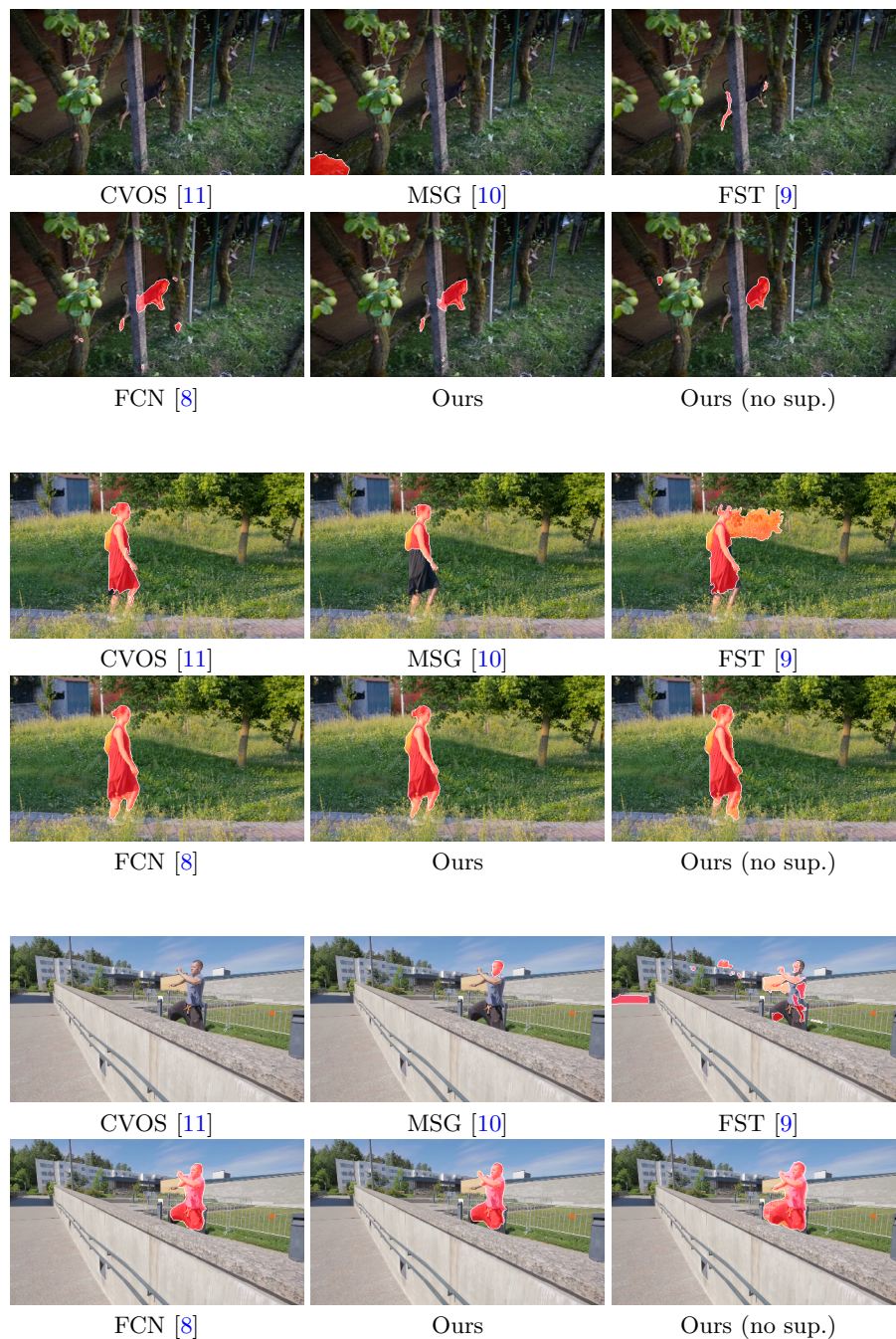


**Fig. 4.** Sample results on the DAVIS dataset with categories shared in the PASCAL VOC dataset. We show that our results with and without weak supervisions have more complete object segments with details.



**Fig. 5.** Sample results on the DAVIS dataset with categories shared in the PASCAL VOC dataset. We show that our results with and without weak supervisions have more complete object segments with details.





**Fig. 6.** Sample results on the DAVIS dataset with categories shared in the PASCAL VOC dataset. We show that our results with and without weak supervisions have more complete object segments with details.

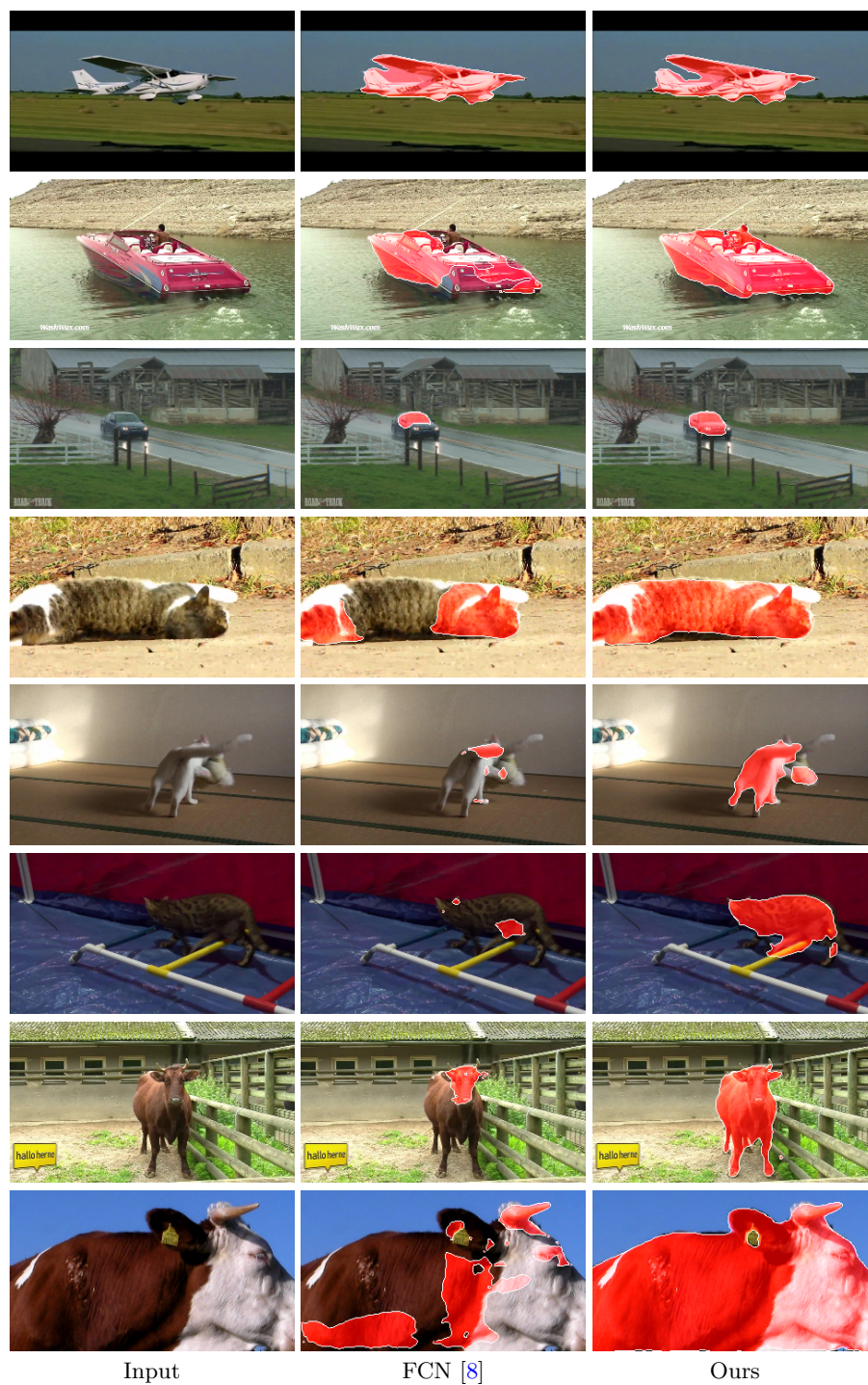
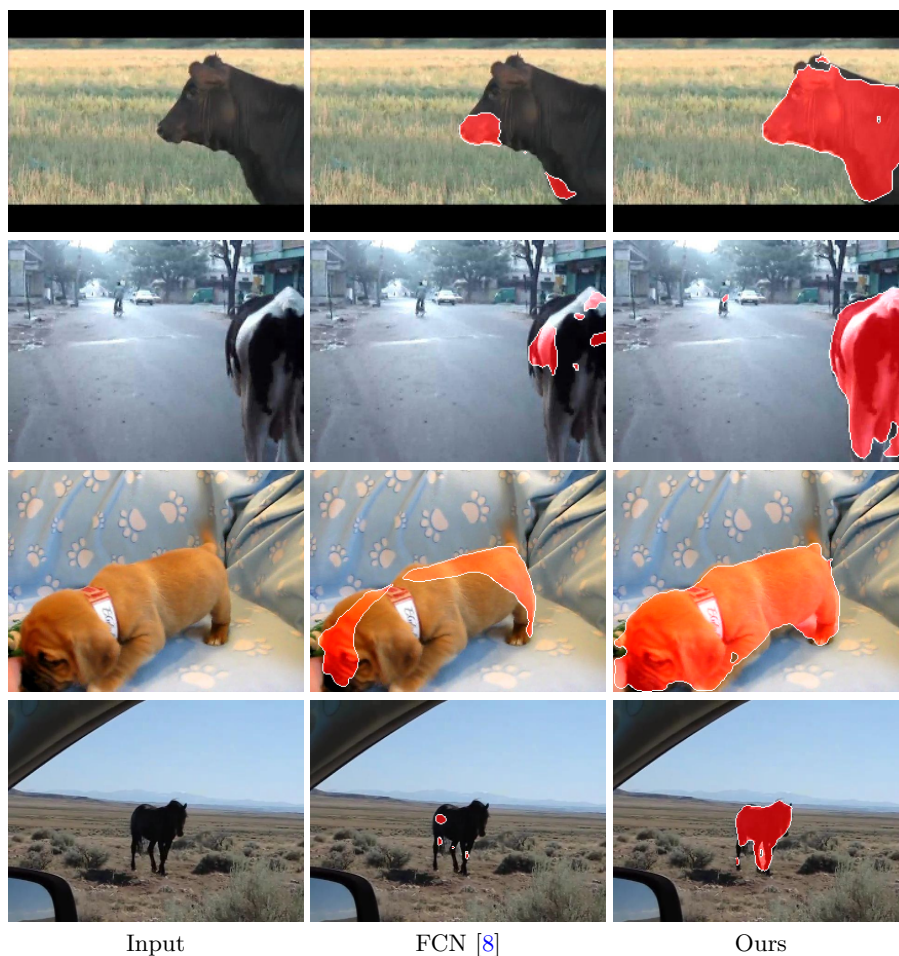
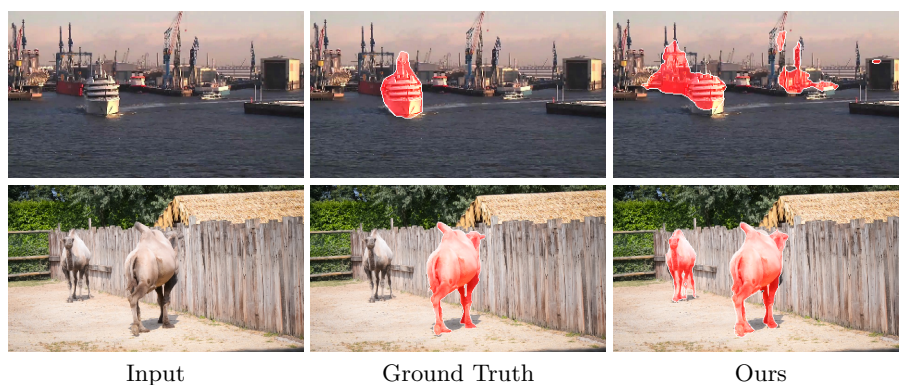


Fig. 7. Sample results on the YouTube-Objects dataset.



Input FCN [8] Ours

**Fig. 8.** Sample results on the YouTube-Objects dataset.



Input Ground Truth Ours

**Fig. 9.** Sample failure cases. Although our results differ from the ground truths, the segmented areas belong to the same semantic category.