

1 We thank the reviewers for valuable feedback. Before addressing individual comments, we clarify common concerns.

2 **Evaluation Protocol:** The most ambitious aim of self-supervised learning is to create universal visual representations
3 that do not require end-to-end fine-tuning, but are instead portable to new tasks merely through additional of a shallow
4 auxiliary network. Moreover, “image-level” vs “pixel-level” training has no bearing on the validity of evaluating with
5 respect to this stricter protocol. Any method that uses a CNN learns more than just “image-level” representations; for
6 example, a CNN’s internal activations can always be converted into per-pixel embeddings using hypercolumns (as we
7 do; see L189). While we agree that adding more fine-tuning results is important for expanding comparison to past work,
8 and will do so, we believe our more ambitious evaluation methodology is fundamental to driving progress in the field.

9 **Dataset Properties:** Progress in supervised object segmentation has been driven by a series of datasets with increasing
10 label complexity per scene, from PASCAL to COCO to LVIS¹. It is reasonable to suspect a similar trend will hold for
11 self-supervised learning, in which case it will be essential to move away from methods such as [19,5,7], which build
12 pipelines dependent upon the single-object bias of the dataset (ImageNet). Corroborating this view are recent results²
13 suggesting even supervised ImageNet pretraining is of minimal value for learning visual tasks on complex scenes.

14 **Experiments:** As we are interested in learned embeddings that *universally* and *directly* fit a broad range of tasks
15 (e.g. instance tracking, segment search, semantic segmentation), we freeze the backbone in most experiments. We have
16 now also run the end-to-end fine-tuning experiment on COCO+VOC pretrained backbone, following the design in [19]
17 to replace DeepLabV3 head with two stacked convolution layers. Results are: ours 47.2 vs MoCo 46.9 mIOU. We use
18 a simpler siamese training architecture, instead of a momentum encoder and memory bank, yet achieve comparable
19 results to MoCo. With frozen pretrained backbone, we also run instance segmentation tasks on COCO2014, following
20 Mask RCNN in [19]. Results of ours vs MoCo: box AP 18.62 vs 14.75 and mask AP 17.82 vs 14.31.

21 Suggested by **R4**, we retrain our model on COCO+VOC with HED edges and achieve 49.9 mIOU in above mentioned
22 end-to-end fine-tuning settings and 48.8 (raised from 46.5) mIOU in the frozen backbone experiment. Suggested by **R3**,
23 to demonstrate the efficiency of proposed pixel-wise optimization, we train MoCo on VOC only for the same number of
24 iterations (32K) and get 26.7 mIOU (ours is 43.5 mIOU) evaluated by freezing the backbone.

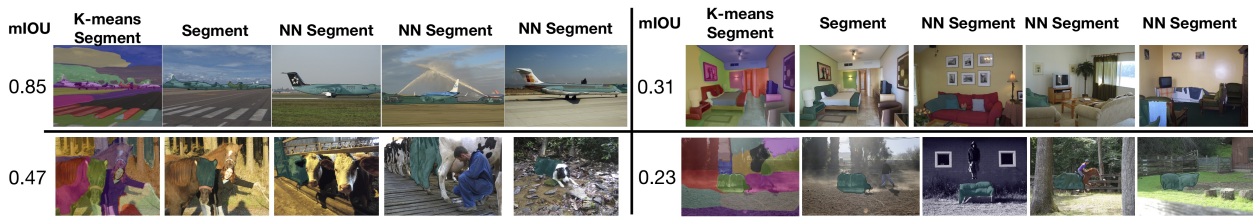


Figure 1: Visualization of success and failure cases denoted by mIOU on segment search task.

25 **R1 – COCO and VOC less curated?** LVIS¹ analyzes statistics of COCO images, finding many instances per-image
26 and a long-tailed category distribution, which supports our characterization. While Instagram-1B does not appear to be
27 publicly available, we worry that Instagram images might suffer from significant photographer bias.

28 **R2, R3, R4 – SegSort.** SegSort focuses on a different task: using contours to refine ImageNet pre-trained features,
29 producing semantic segmentation. Our task is to learn pixel-wise semantic-aware embeddings from scratch. To compare,
30 we evaluate SegSort with its backbone initialized from scratch. Performance in “semantic segment retrieval”, the truly
31 unsupervised equivalent of SegSort’s “unsupervised semantic segmentation” is this lower than their reported scores.

32 **R1, R2 – Partial ImageNet results.** We will update the final version to reflect the full 200 training epochs.

33 **R1, R2 – Sampling details.** We first sample regions, then a fixed number of pixels within chosen regions. Probability
34 should be summed over regions, i.e. $\sum_{R_b} P(j \in Pos(i)) = 1$. Pixels could be sampled as positive and negative at the
35 same time. Instead of a memory bank, we compute the representation from both original and augmented images.

36 **R1 – Distance in tree.** We compute $d_T(R_i, R_j)$ following Sec. 4.2 of [2], originally introduced in Sec. 2-4 of [1].

37 **R3 – Hyperparameters.** We pick σ_p using limited search (0.2 - 1.0, spaced by 0.2). For MoCo training on COCO +
38 VOC, we also ran a limited hyperparameter search and did not see significant impact on results.

39 **Citations. R1:** We will add suggested papers. **R3:** Missed citation is indeed Li *et al.*, CVPR 2016.

40 **Misc. R1:** Table 1 row 9 should be 36.15. We will update Fig. 3, 4, 5 with best/worst cases (for Fig. 4, some examples
41 are shown above in Fig. 1). **R2:** ‘Image’ under ‘Cross-Image Negative Sample’ refers to treating the image-based
42 embedding as negative features. ‘spatially-extended’ means we preserve the spatial dimension in the embedding and
43 output a feature map. **R3:** We fine-tune on PASCAL *train_aug* and evaluate on PASCAL *val*.

¹A. Gupta, P. Dollar, R. Girshick. LVIS: A Dataset for Large Vocabulary Instance Segmentation. CVPR, 2019.

²K. He, R. Girshick, P. Dollar. Rethinking ImageNet Pre-training. ICCV, 2019.