

1 We thank the reviewers for their feedback. We appreciate that all the reviewers agree on the value of our analysis
 2 and they like the intuitive idea of augmentation via videos. We also acknowledge that more experiments in 4.2 (R3)
 3 would help with justification/interpretation of results (R2). We are pleased to report that we have conducted several
 4 experiments that address *all* the concerns of the reviewers. Specifically, we perform two additional experiments to
 5 highlight how aggressive augmentation hurts (and it is not just the domain-gap). We have also conducted experiments
 6 to show that on a different downstream task, *video-based augmentation outperforms even ImageNet-based MoCo*.
 7 The goal of this paper is to present a critical analysis so that the community can also introspect as we make rapid
 8 progress on the topic of representation learning. We hope the discussion here and additional experiments will convince
 9 the reviewers and AC that this message deserves a wider audience at NeurIPS.

10 **R2: Domain Gap or Augmentation effect? R3: Additional Experiment in 4.2** We agree: there is some ambiguity
 11 since the training data for COCO-cropped model is domain-aligned with PASCAL-cropped test data (unlike COCO-
 12 full). The drop in performance could be attributed to domain gap between COCO-full and PASCAL-cropped. So, we
 13 conducted a new experiment that cannot be explained by domain-gap.

14 **Setup:** Consider the subset of Pascal VOC07 images which depict either table or chair in the image, but not both (i.e
 15 scene-centric images containing only one of the frequently co-occurring pair of objects).

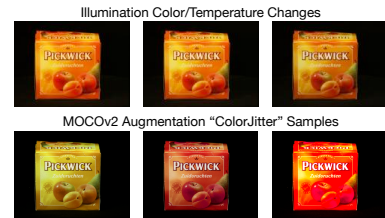
16 **Result:** On the table vs chair full image classification task, the representation trained on COCO-Boxes outperforms
 17 full COCO-image pre-training - 74.92 vs 73.64 mAP . Note that the test domain (PASCAL-full images) is aligned
 18 with training domain (COCO-full images). Yet, COCO-Boxes model outperforms COCO-full. This indicates that the
 19 problem lies with aligned representations of co-occurring objects (e.g., chairs and tables) as explained in the paper.

20 To support it even further, we do another experiment. We compare super-
 21 vised learning on ImageNet with MOCO on ImageNet but for the task of
 22 PASCAL-part classification. Note when training on ImageNet, MOCO’s
 23 cropping augmentations will learn embeddings that have similar represen-
 24 tation for different parts. Hence, our augmentation interpretation predicts

Method	Dog Parts	Cat Parts
ImageNet R50	69.14	71.31
ImageNet MOCOv2	51.91	56.26
(Ours) Region Tracker	54.70	58.81

25 MOCO should not be as robust as supervised models on part-classification task. And indeed! the results indicate MOCO
 26 performance falls dramatically on part-classification.

27 **PIRL/MOCO should have higher Illumination Color Invariance (R2)** The
 28 natural illumination color changes for measuring invariances are restricted to
 29 illumination *temperature* changes from 2175K to 3075K - this is also naturally
 30 captured in the ImageNet dataset. The synthetically color-jittered samples
 31 are significantly different from these images depicting arbitrary color changes
 32 (see adjoining figure). On the other hand, occlusion is very easy to accurately
 33 synthesize by simply cropping images. And therefore the cropping augmentation strategy indeed leads to high occlusion
 34 invariance even on natural images.



35 **Video-based Approach - Aggressive Augmentation and Performance (R1, R2)** Our method only uses the
 36 aggressive augmentation in the Frame Temporal Invariance loss. The other component of the loss involving region
 37 tracks does not employ the aggressive cropping strategy. Including the aggressive augmentations ensures that we can
 38 take advantage of the occlusion invariance that it induces, while being constrained by the region-level loss. This extra
 39 constraint ensures the best of both worlds. This method was presented as a proof-of-concept to demonstrate that the
 40 appropriate invariances can be induced by using video data. Therefore, we did not scale up the data for this approach
 41 up to be comparable to ImageNet in terms of dataset size. However, we are pleased to report (as R2 points out) that
 42 on a downstream task which is ill-suited to ImageNet-MOCO (due to aggressive augmentation), our video approach
 43 outperforms it (even with order of magnitude less data). The results on VOC-part classification shows that trend (See
 44 table above), also demonstrating that our approach is not affected by the aggressive augmentation.

45 **RIS for Measuring Invariances (R3)** Goodfellow et. al. [26] explains how the firing frequency can be used to
 46 measure invariances. We directly adopt the same principle. To the best of our knowledge, this is the most relevant
 47 metric proposed in past literature. It is true that two somewhat distinct representations after thresholding can have
 48 the same invariance under this metric. However, this is an intended feature proposed in [26] since most downstream
 49 classifiers would threshold a function of the neurons (generally linear classifiers and non-linear models in some cases).

50 **Object Detectors as Objectness Prior (R1):** We use a fully-unsupervised class-agnostic region proposal method -
 51 Selective Search. One should view this as a simple low-level preprocessing operation which ensures object-centric
 52 signals necessary for the augmentations. Since it uses no learning, we believe the comparisons are fair.

53 **Argument on Capacity (R1):** In light of new experiments, we agree with R1 and plan to remove this as it has no
 54 connection to main story. **Finetuning representations (R4):** For this work, we wished to analyze the representations
 55 learned in the contrastive learning framework. Recent SSL benchmarks have suggested avoiding finetuning [25] since it
 56 would cause the representations to deviate from the initially learned representations, leading to spurious inferences.