

1 Before addressing the issues raised by the reviewers, we present the
 2 results of using our self-supervised representation for segmentation mask
 3 propagation on the DAVIS dataset (Table 1), in comparison to a SOTA
 4 self-supervised method [49] and the ImageNet pre-trained representation.

5 We sincerely urge the area chair and reviewers to evaluate this paper in
 6 the context of advances accomplished by our proposed self-supervised
 7 learning method. This is the **first** paper to show large improvements over
 8 an ImageNet pre-trained representation. The contributions are significant, since the performance of the existing SOTA
 9 [49] is far below our method and also below that of the ImageNet pre-training method.

10 **Novelty and Contributions. (R1, R2)** We would like to remind the reviewers and the area chair of the challenges
 11 and difficulties in leveraging various self-supervision signals. Most existing methods still focus solely on applying
 12 multiple losses, e.g., [17, 43, 44, 49], and achieve marginal improvements. Instead of performing multi-loss training,
 13 we propose a novel framework, in which we build one task upon another that progressively improve each other during
 14 training. The shared affinity matrix bridges these tasks, and facilitates iterative improvements. The proposed framework
 15 not only improves previous self-supervised approaches significantly, but outperforms supervised learning with human
 16 supervision by a large margin. These contributions are significant in the field of self-supervised learning.

17 The contributions of this work are also demonstrated by our ablation study, i.e., Table 2 in the paper. All the proposed
 18 components, e.g., coarse localization, concentration and orthogonal regularization, contribute to the performance gain.
 19 We note that these components are novel and have not been explored in prior work.

20 **Similarity to [49] (R2)** The key ideas of our paper are not similar to [49]: (i) [49] matches patches instead of pixels,
 21 and the entire fine-grained matching part is missing; (ii) [49] models the locations independently via a STN, while we
 22 track both features and locations uniformly via the same affinity matrix; (iii) The cycle-consistency in our method refers
 23 to the orthogonal regularization of the affinity matrix, which is very much different from the cycle-consistency tracking
 24 loss in [49]. The achieved performance gain of this work comes from all the above-mentioned algorithmic components
 25 and concentration regularization, rather than engineering work.

Table 1: DAVIS-2017 segmentation results.

Model	\mathcal{J} (Mean)
Self-supervised, SOTA [49]	43.0
ImageNet Representation	49.4
Self-supervised, Ours	57.7

Table 2: Comparison against optical flow methods. Table 3: Ablation study on temperature (no track in testing).

Model	\mathcal{J} (Mean)	\mathcal{F} (Mean)
FlowNet2 [16]	26.7	25.2
PWC-Net [39]	35.2	37.4
Ours	57.7	60.0

Temperature	\mathcal{J} (Mean)	\mathcal{F} (Mean)
1	56.8	59.5
2	52.3	56.2
10	51.0	55.5

26 **Which methods should the work compare with? (R2)** We note that the focus of this work is not on learning
 27 unsupervised flow. Rather, it aims to achieve the same goal as presented in [43] and [49] – learning unsupervised
 28 correspondences. These two tasks are significantly different: one focuses on learning a feature matching network
 29 that can track regions and pixels between frames over a long period, while the other emphasizes modeling subtle
 30 displacements between adjacent frames. Taking Table 2 as an example, for the segmentation propagation task on the
 31 DAVIS dataset, even the SOTA **fully-supervised** flow methods perform much worse than any of the unsupervised
 32 correspondence methods. All the latest methods related to unsupervised video correspondences are evaluated, including
 33 several concurrent works, e.g., [7,8] in the supplementary material. The comparisons with UDT [44] are thoroughly
 34 discussed in Section 4.4.

35 In the following, we address the other comments by reviewers.

36 **Sharpness in the Softmax layer. (R1)** We experimented with various temperatures in the softmax function and
 37 found that setting it to 1 achieves the best results. See Table 3.

38 **Equation (6) makes loss non-smooth and non-differentiable. (R1)** With Eq. (6), we only assign a penalty to pixels
 39 that move outside the bounding box, which is achieved by multiplying a rectangular mask to the loss term. In practice,
 40 there is no need to specify gradients on the boundary to enforce smoothness. We will clarify this in the revised paper.

41 **Scale estimation in equation 4. (R3)** We compute the scale strictly according to Eq. (4), rather than computing the
 42 maximum distance. Please refer to Figure 3 in the supplementary for visualizations of the estimated scales.

43 **Which dataset to train the auto-encoder? (R3)** The MSCOCO dataset is large and diverse for training an image
 44 auto-encoder, e.g., with a shallow 6-layer encoder. The effectiveness of our trained auto-encoder is also validated in
 45 Table 2 (g). Under the same settings, we achieve 11.1% higher accuracy than [43] that does not use the auto-encoder.

46 **Evaluating learned representations on the image based tasks (R3)** Transferring the learned representations from
 47 video correspondence to image based task is certainly interesting. We will explore it in the future.