
Self-Supervised Spatial-Temporal Feature Learning for Video Correspondence

— NeurIPS 2022 Supplementary Material

Anonymous Author(s)

Affiliation

Address

email

1 The supplementary material contains: 1) ablation study of different contrastive models; 2) ablation
2 study of entropy-based selection; 3) more qualitative examples for video object segmentation.

3 1 Ablation study of different contrastive models

4 Given a query point randomly sampled in the target frame, we visualize the result of computing the
5 local correlation and global correlation map *w.r.t.* reference frame. The dashed line in red represents
6 the range of computing correlation map *w.r.t.* query point. The reference frame is randomly sampled
7 in the memory bank of inference strategy. Given a query point randomly sampled in the target frame,
8 we visualize the result of computing the local correlation and global correlation map *w.r.t.* reference
9 frame. The dashed line in red represents the range of computing correlation map *w.r.t.* query point.
10 The reference frame is randomly sampled in the memory bank of inference strategy. Given a query
11 point randomly sampled in the target frame, we visualize the result of computing the local correlation
12 and global correlation map *w.r.t.* reference frame. The dashed line in red represents the range of
13 computing correlation map *w.r.t.* query point. The reference frame is randomly sampled in the
14 memory bank of inference strategy

15 2 Ablation study of entropy-based selection

16 We make detailed ablation study for our entropy-
17 based selection in terms of both visual perception
18 and quantitative comparison. As shown in Figure 1,
19 the entropy map has a higher response on moving
20 objects involved in severe deformation and occlu-
21 sions, which should be paid more attention to. In
22 Table 1, we adopt different thresholds to generate
23 the mask m with high entropy. The baseline is to
24 apply local correlation distillation for all queries,
25 *i.e.*, $T = 1.0$. When setting T with 0.1, the perfor-
26 mance drops to 67.4% due to the underutilization of
27 the supervision from finest pyramid level. The results of setting T with 0.4 and 0.7 indicate applying
28 distillation in the region with high entropy exhibits a performance gain.

Method	Dataset	$\mathcal{J} \& \mathcal{F}_m \uparrow$
$T = 0.1$	YTV	67.4
$T = 0.4$	YTV	68.3
$T = 0.7$	YTV	69.0
$T = 1.0$	YTV	68.1

Table 1: **The quantitative results on the validation set of DAVIS-2017.** The Dataset represents dataset(s) used for training. YTV: YouTube-VOS []

29 3 More qualitative examples for video object segmentation

30 We do the most of experiments

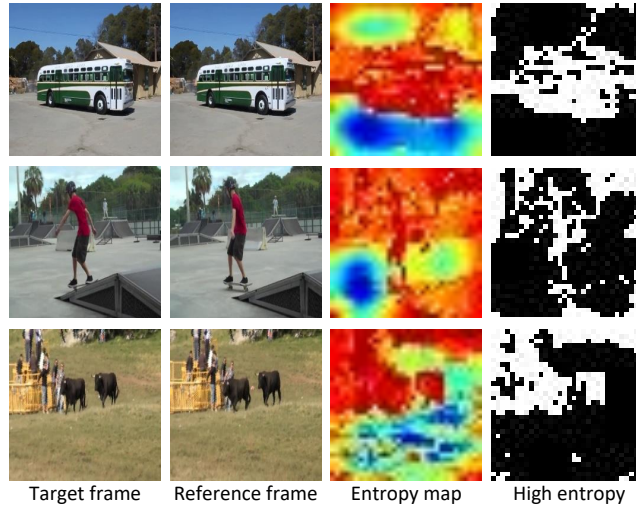


Figure 1: **Visualization of the entropy map.** We compute the entropy for each query in target frame using Eq 8 in our main paper. The mask with high entropy is generated by setting a threshold.