

Contrastive Transformation for Self-supervised Correspondence Learning

Ning Wang¹, Wengang Zhou^{1,2}, Houqiang Li^{1,2}

¹ CAS Key Laboratory of GIPAS, University of Science and Technology of China

² Institute of Artificial Intelligence, Hefei Comprehensive National Science Center
wn6149@mail.usstc.edu.cn, {zhwg, lihq}@ustc.edu.cn

Abstract

In this paper, we focus on the self-supervised learning of visual correspondence using unlabeled videos in the wild. Our method simultaneously considers intra- and inter-video representation associations for reliable correspondence estimation. The intra-video learning transforms the image contents across frames within a single video via the frame pair-wise affinity. To obtain the discriminative representation for instance-level separation, we go beyond the intra-video analysis and construct the inter-video affinity to facilitate the contrastive transformation across different videos. By forcing the transformation consistency between intra- and inter-video levels, the fine-grained correspondence associations are well preserved and the instance-level feature discrimination is effectively reinforced. Our simple framework outperforms the recent self-supervised correspondence methods on a range of visual tasks including video object tracking (VOT), video object segmentation (VOS), pose keypoint tracking, *etc.* It is worth mentioning that our method also surpasses the fully-supervised affinity representation (*e.g.*, ResNet) and performs competitively against the recent fully-supervised algorithms designed for the specific tasks (*e.g.*, VOT and VOS).

1 Introduction

Learning representations for visual correspondence is a long-standing problem in computer vision, which is closely related to many vision tasks including video object tracking, keypoint tracking, and optical flow estimation, *etc.* This task is challenging due to the factors such as viewpoint change, distractors, and background clutter.

Correspondence estimation generally requires human annotations for model training. Collecting dense annotations, especially for large-scale datasets, requires costly human efforts. To leverage the large volume of raw videos in the wild, the recent advances focus on self-supervised correspondence learning by exploring the inherent relationships within the unlabeled videos. In (Wang, Jabri, and Efros 2019), the temporal cycle-consistency is utilized to self-supervise the feature representation learning. To be specific, the correct patch-level or pixel-wise associations between two successive frames should match bi-directionally

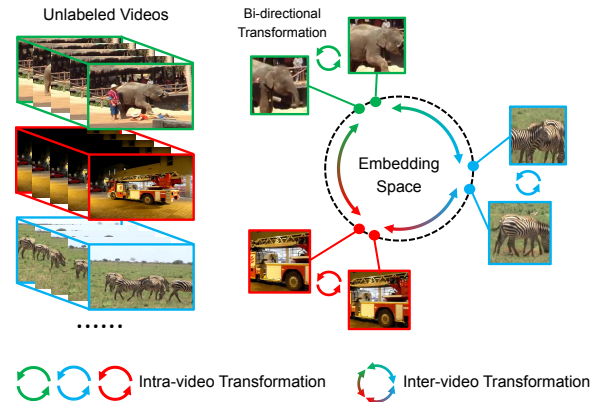


Figure 1: The proposed approach targets at learning correspondence using unlabeled videos. Previous works mainly focus on the content transformation within each video clip. Our framework simultaneously tracks (intra-video level) and spreads (inter-video level) the feature embeddings to preserve the fine-grained matching capability while encouraging the contrastive embedding learning.

in both forward and backward tracking trajectories. The bi-directional matching is realized via a frame-level affinity matrix, which represents the pixel pair-wise similarity between two frames. In (Vondrick et al. 2018; Li et al. 2019), this affinity is also utilized to achieve the content transformation between two frames for self-supervision. A straightforward transformation within videos is the color/RGB information. More specifically, the pixel colors in a target frame can be “copied” (or transformed) from the pixels in a reference frame. By minimizing the differences between the transformed and the true colors of the target frame, the backbone network is forced to learn robust feature embeddings for identifying correspondence across frames in a self-supervised manner.

In spite of the impressive performance, existing unsupervised correspondence algorithms put all the emphasis on the intra-video analysis. Since the scenario in one video is generally stable and changeless, establishing the correspondence within the same videos is less challenging and inevitably hinders the discrimination potential of learned fea-

ture embeddings. In this work, we go beyond the intra-video correspondence learning by further considering the inter-video level embedding separation of different instance objects. Our method is largely inspired by the recent success of contrastive learning (He et al. 2020; Chen et al. 2020), which aims at minimizing the agreement between different augmented versions of the same image via a contrastive loss (Hadsell, Chopra, and Lecun 2006). Nevertheless, there are two obvious gaps between contrastive learning and correspondence learning. First, classic contrastive learning relies on the augmented still images, but how to adapt it to the video-level correspondence scenario is rarely explored. Second, their optimization goals are somewhat conflicting. Contrastive learning targets at positive concentration and negative separation, ignoring the pixel-to-pixel relevance among the positive embeddings. In contrast, correspondence learning aims at identifying fine-grained matching.

In this work, we aim to narrow the above domain gaps by absorbing the core contrastive ideas for correspondence estimation. To transfer the contrastive learning from the image domain to the video domain, we leverage the patch-level tracking to acquire matched image pairs in unlabeled videos. Consequently, our method captures the real target appearance changes reside in the video sequences without augmenting the still images using empirical rules (*e.g.*, scaling and rotation). Furthermore, we propose the inter-video transformation, which is consistent with the correspondence learning in terms of the optimization goal while preserving the contrastive characteristic among different instance embeddings. In our framework, similar to previous arts (Vondrick et al. 2018; Li et al. 2019), the image pixels should match their counterpart pixels in the current video to satisfy the self-supervision. Besides, these pixels are also forced to mismatch the pixels in other videos to reinforce the instance-level discrimination, which is formulated in the contrastive transformation across a batch of videos, as shown in Figure 1. By virtue of the intra-inter transformation consistency as well as the sparsity constraint for the inter-video affinity, our framework encourages the contrastive embedding learning within the correspondence framework.

In summary, the main contribution of this work lies in the contrastive framework for self-supervised correspondence learning. 1) By joint unsupervised tracking and contrastive transformation, our approach extends the classic contrastive idea to the temporal domain. 2) To bridge the domain gap between two diverse tasks, we propose the intra-inter transformation consistency, which differs from contrastive learning but absorbs its core motivation for correspondence tasks. 3) Last but not least, we verify the proposed approach in a series of correspondence-related tasks including video object segmentation, pose tracking, object tracking, *etc.* Our approach consistently outperforms previous state-of-the-art self-supervised approaches and is even comparable with some task-specific fully-supervised algorithms.

2 Related Work

In this section, we briefly review the related methods including unsupervised representation learning, self-supervised correspondence learning, and contrastive learning.

Unsupervised Representation Learning. Learning representations from unlabeled images or videos has been widely studied. Unsupervised approaches explore the inherent information inside images or videos as the supervisory signals from different perspectives, such as frame sorting (Lee et al. 2017), image content recovering (Pathak et al. 2016), deep clustering (Caron et al. 2018), affinity diffusion (Huang et al. 2020), motion modeling (Pathak et al. 2017; Tung et al. 2017), and bi-directional flow estimation (Meister, Hur, and Roth 2018). These methods learn an unsupervised feature extractor, which can be generalized to different tasks by further fine-tuning using a small set of labeled samples. In this work, we focus on a sub-area in the unsupervised family, *i.e.*, learning features for fine-grained pixel matching without task-specific fine-tuning. Our framework shares partial insight with (Wang and Gupta 2015), which utilizes off-the-shelf visual trackers for data pre-processing. Differently, we jointly track and spread feature embeddings in an end-to-end manner for complementary learning. Our method is also motivated by the contrastive learning (Den Oord, Li, and Vinyals 2018), another popular framework in the unsupervised learning family. In the following, we will detailedly discuss correspondence learning and contrastive learning.

Self-supervised Correspondence Learning. Learning temporal correspondence is widely explored in the visual object tracking (VOT), video object segmentation (VOS), and flow estimation (Dosovitskiy et al. 2015) tasks. VOT aims to locate the target box in each frame based on the initial target box, while VOS propagates the initial target mask. To avoid expensive manual annotations, self-supervised approaches have attracted increasing attention. In (Vondrick et al. 2018), based on the frame-wise affinity, the pixel colors from the reference frame are transferred to the target frame as self-supervisory signals. Wang *et al.* (Wang, Jabri, and Efros 2019) conduct the forward-backward tracking in unlabeled videos and leverage the inconsistency between the start and end points to optimize the feature representation. UDT algorithm (Wang et al. 2019) leverages a similar bi-directional tracking idea and composes the correlation filter for unsupervised tracker training. In (Yang, Zhang, and Zhang 2019), an unsupervised tracker is trained via incremental learning using a single movie. Recently, Li *et al.* (Li et al. 2019) combine the object-level and fine-grained correspondence in a coarse-to-fine fashion and shows notable performance improvements. In (Jabri, Owens, and Efros 2020), space-time correspondence learning is formulated as a contrastive random walk and shows impressive results. Despite the success of the above methods, they put the main emphasis on the intra-video self-supervision. Our approach takes a step further by simultaneously exploiting the intra-video and inter-video consistency to learn more discriminative feature embeddings. Therefore, previous intra-video based approaches can be regarded as one part of our framework.

Contrastive Learning. Contrastive learning is a popular unsupervised learning paradigm, which aims to enlarge the embedding disagreements of different instances for representation learning (Den Oord, Li, and Vinyals 2018; Ye et al. 2019; Hjelm et al. 2019). Based on the contrastive framework, the recent SimCLR method (Chen et al. 2020) signifi-

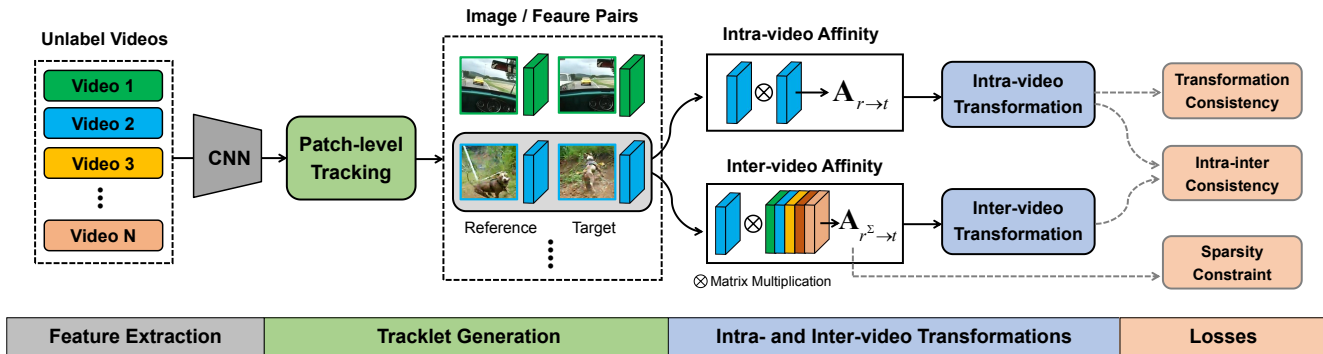


Figure 2: An overview of the proposed framework. Given a batch of videos, we first do patch-level tracking to generate image pairs. Then, intra- and inter-video transformations are conducted for each video in the mini-batch. Finally, except the intra-video self-supervision, we introduce the intra-inter consistency and sparsity constraint to reinforce the embedding discrimination.

cantly narrows the performance gap between supervised and unsupervised models. He *et al.* (He et al. 2020) propose the MoCo algorithm to fully exploit the negative samples in the memory bank. Inspired by the recent success of contrastive learning, we also involve plentiful negative samples for discriminative feature learning. Compared with existing contrastive methods, one major difference is our method jointly tracks and spreads feature embeddings in the video domain. Therefore, our method captures the temporally changed appearance variations instead of manually augmenting the still images. Besides, instead of using a standard contrastive loss (Hadsell, Chopra, and Lecun 2006), we incorporate the contrastive idea into the correspondence task by a conceptually simple yet effective contrastive transformation mechanism to narrow the domain gap.

3 Methodology

An overview of our framework is shown in Figure 2. Given a batch of videos, we first crop the adjacent image patches via patch-level tracking, which ensures the image pairs have similar contents and facilitates the later transformations. For each image pair, we consider the intra-video bi-directional transformation. Furthermore, we introduce irrelevant images from other videos to conduct the inter-video transformation for contrastive embedding learning. The final training objectives include the intra-video self-supervision, intra-inter transformation consistency, and sparsity regularization for the batch-level affinity.

3.1 Revisiting Affinity-based Transformation

Given a pair of video frames, the pixel colors (*e.g.*, RGB values) in one frame can be copied from the pixels from another frame. This is based on the assumption that the contents in two successive video frames are coherent. The above frame reconstruction (pixel copy) operation can be expressed via a linear transformation with the affinity matrix $\mathbf{A}_{r \rightarrow t}$, which describes the copy process from a reference frame to a target frame (Vondrick et al. 2018; Liu et al. 2018).

A general option for the similarity measurement in the affinity matrix is the dot product between feature embeddings. In this work, we follow previous arts (Vondrick et al.

2018; Wang, Jabri, and Efros 2019; Li et al. 2019) to construct the following affinity matrix:

$$\mathbf{A}_{r \rightarrow t}(i, j) = \frac{\exp(\mathbf{f}_t(i)^\top \mathbf{f}_r(j))}{\sum_j \exp(\mathbf{f}_t(i)^\top \mathbf{f}_r(j))}, \quad (1)$$

where $\mathbf{f}_t \in \mathbb{R}^{C \times N_1}$ and $\mathbf{f}_r \in \mathbb{R}^{C \times N_2}$ denote flattened feature maps with C channels of target and reference frames, respectively. With the spatial index $i \in [1, N_1]$ and $j \in [1, N_2]$, $\mathbf{A}_{r \rightarrow t} \in \mathbb{R}^{N_1 \times N_2}$ is normalized by the softmax over the spatial dimension of \mathbf{f}_r .

Leveraging the above affinity, we can freely transform various information from the reference frame to the target frame by $\hat{\mathbf{L}}_t = \mathbf{A}_{r \rightarrow t} \mathbf{L}_r$, where \mathbf{L}_r can be any associated labels of the reference frame (*e.g.*, semantic mask, pixel color, and pixel location). Since we naturally know the color information of the target frame, one free self-supervisory signal is color (Vondrick et al. 2018). The goal of such an affinity-based transformation framework is to train a good feature extractor for affinity computation.

3.2 Contrastive Pair Generation

A vital step in contrastive frameworks is building positive image pairs via data augmentation. We free this necessity by exploring the temporal content consistency resides in the videos. To this end, for each video, we first utilize the patch-level tracking to acquire a pair of high-quality image patches with similar content. Based on the matched pairs, we then conduct the contrastive transformation.

Given a randomly cropped patch in the reference frame, we aim to localize the best matched patch in the target frame, as shown in Figure 2. Similar to Eq. 1, we compute a patch-to-frame affinity between the features of a random patch in the reference frame and the features of the whole target frame. Based on this affinity, in the target frame, we can identify some target pixels most similar to the reference pixels, and average these pixel coordinates as the tracked target center. We also estimate the patch scale variation following UVC approach (Li et al. 2019). Then we crop this patch and combine it with the reference patch to form an image pair.

3.3 Intra- and Inter-video Transformations

Intra-video. After obtaining a pair of matched feature maps via patch-level tracking, we compute their fined-grained affinity $\mathbf{A}_{r \rightarrow t}$ according to Eq. 1. Based on this intra-video affinity, we can easily transform the image contents from the reference patch to the target patch within a single video clip.

Inter-video. The key success of the aforementioned affinity-based transformation lies in the embedding discrimination among plentiful subpixels to achieve the accurate label copy. Nevertheless, within a pair of small patch regions, the image contents are highly correlated and even only cover a subregion of a large object, struggling to contain diverse visual patterns. The rarely existing negative pixels from other instance objects heavily hinder the embedding learning.

In the following, we improve the existing framework by introducing another inter-video transformation to achieve the contrastive embedding learning. The inter-video affinity is defined as follows:

$$\mathbf{A}_{r \rightarrow t}(i, j) = \frac{\exp(\mathbf{f}_t(i)^\top \mathbf{f}_r^\Sigma(j))}{\sum_j \exp(\mathbf{f}_t(i)^\top \mathbf{f}_r^\Sigma(j))}, \quad (2)$$

where \mathbf{f}_r^Σ is the concatenation of the reference features from different videos in the spatial dimension, *i.e.*, $\mathbf{f}_r^\Sigma = \text{Concat}(\mathbf{f}_r^1, \dots, \mathbf{f}_r^n)$. For a mini-batch with n videos, the spatial index $i \in [1, N_1]$ and $j \in [1, nN_2]$.

Rationale Analysis. Inter-video transformation is an extension of intra-video transformation. By decomposing the reference feature embeddings $\mathbf{f}_r^\Sigma \in \mathbb{R}^{C \times nN_2}$ into positive and negative, \mathbf{f}_r^Σ can be expressed as $\mathbf{f}_r^\Sigma = \text{Concat}(\mathbf{f}_r^+, \mathbf{f}_r^-)$, where $\mathbf{f}_r^+ \in \mathbb{R}^{C \times N_2}$ denotes the only positive reference feature related to the target frame feature while $\mathbf{f}_r^- \in \mathbb{R}^{C \times (n-1)N_2}$ is the concatenation of negative ones from unrelated videos in the mini-batch. As a result, the computed affinity $\mathbf{A}_{r \rightarrow t} \in \mathbb{R}^{N_1 \times nN_2}$ can be regarded as an ensemble of multiple sub-affinities, as shown in Figure 3. Our goal is to build such a batch-level affinity for discriminative representation learning.

To facilitate the later descriptions, we also divide the inter-video affinity $\mathbf{A}_{r \rightarrow t}$ as a combination of positive and negative sub-affinities:

$$\mathbf{A}_{r \rightarrow t} = \text{Concat}(\mathbf{A}_{r^+ \rightarrow t}, \mathbf{A}_{r^- \rightarrow t}), \quad (3)$$

where $\mathbf{A}_{r^+ \rightarrow t} \in \mathbb{R}^{N_1 \times N_2}$ and $\mathbf{A}_{r^- \rightarrow t} \in \mathbb{R}^{N_1 \times (n-1)N_2}$ are the positive and negative sub-affinities, respectively. Ideally, sub-affinity $\mathbf{A}_{r^+ \rightarrow t}$ should be close to the intra-video affinity and $\mathbf{A}_{r^- \rightarrow t}$ is expected to be a zero-like matrix. Nevertheless, with the inclusion of noisy reference features \mathbf{f}_r^- , the positive sub-affinity $\mathbf{A}_{r^+ \rightarrow t}$ inevitably degenerates in comparison with the intra-video affinity $\mathbf{A}_{r \rightarrow t}$, as shown in Figure 3. In the following, we present the intra-inter transformation consistency to encourage contrastive embedding learning within the correspondence learning task.

3.4 Training Objectives

To achieve the high-quality frame reconstruction, following (Li et al. 2019), we pre-train an encoder and a decoder using

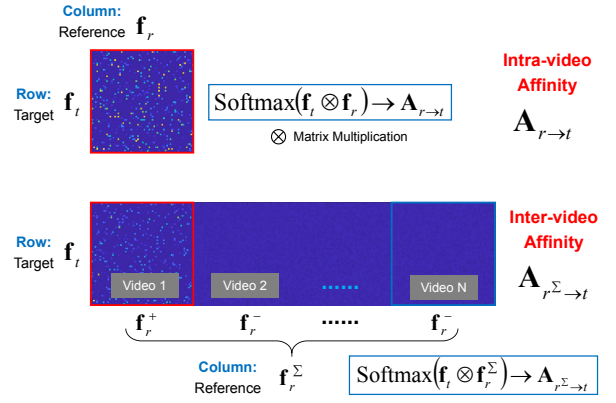


Figure 3: Comparison between intra-video affinity (top) and inter-video affinity (bottom). Best view in zoom in.

still images on the COCO dataset (Lin et al. 2014) to perform the feature-level transformation. The pre-trained encoder and decoder networks are frozen without further optimization in our framework. The goal is to train the backbone network for correspondence estimation (*i.e.*, affinity computation). In the following, the encoded features of the reference image \mathbf{I}_r is denoted as $\mathbf{E}_r = \text{Encoder}(\mathbf{I}_r)$.

Intra-video Self-supervision. Leveraging the intra-video affinity $\mathbf{A}_{r \rightarrow t}$ as well as the encoded reference feature \mathbf{E}_r , the transformed target image can be computed via $\hat{\mathbf{I}}_{r \rightarrow t} = \text{Decoder}(\mathbf{A}_{r \rightarrow t} \mathbf{E}_r)$. Ideally, the transformed target frame should be consistent with the original target frame. As a consequence, the intra-video self-supervisory loss is defined as follows:

$$\mathcal{L}_{\text{self}} = \|\hat{\mathbf{I}}_{r \rightarrow t} - \mathbf{I}_t\|_1. \quad (4)$$

Intra-inter Consistency. Leveraging the inter-video affinity $\mathbf{A}_{r \rightarrow t}$ and the encoded reference features \mathbf{E}_r^Σ from a batch of videos, *i.e.*, $\mathbf{E}_r^\Sigma = \text{Concat}(\mathbf{E}_r^1, \dots, \mathbf{E}_r^n)$, the corresponding transformed target image can be computed via $\hat{\mathbf{I}}_{r \rightarrow t} = \text{Decoder}(\mathbf{A}_{r \rightarrow t} \mathbf{E}_r^\Sigma)$. This inter-video transformation is shown in Figure 4. The reference features from other videos are considered as negative embeddings. The learned inter-video affinity is expected to exclude unrelated embeddings for transformation fidelity. Therefore, the transformed images via intra-video affinity and inter-video affinity should be consistent:

$$\mathcal{L}_{\text{intra-inter}} = \|\hat{\mathbf{I}}_{r \rightarrow t} - \hat{\mathbf{I}}_{r^\Sigma \rightarrow t}\|_1. \quad (5)$$

The above loss encourages both positive feature invariance and negative embedding separation.

Sparsity Constraint. To further enlarge the disagreements among different video features, we force the sub-affinity in the inter-video affinity $\mathbf{A}_{r \rightarrow t}$ to be sparse via

$$\mathcal{L}_{\text{sparse}} = \|\mathbf{A}_{r^- \rightarrow t}\|_1, \quad (6)$$

where $\mathbf{A}_{r^- \rightarrow t}$ is the negative sub-affinity in Eq. 3.

Other Regularizations. Following previous works (Li et al. 2019; Wang, Jabri, and Efros 2019), we also utilize the

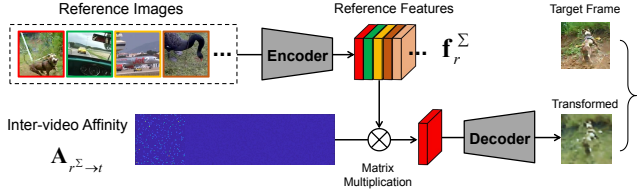


Figure 4: Illustration of the inter-video transformation.

cycle-consistency (bi-directional matching) between two frames, which equals forcing the affinity matrix to be orthogonal, *i.e.*, $\mathbf{A}_{r \rightarrow t}^{-1} = \mathbf{A}_{t \rightarrow r}$. Besides, the concentration regularization proposed in (Li et al. 2019) is also added. These two regularizations are combined and denoted as $\mathcal{L}_{\text{others}}$.

Final Objective. The final training objective is the combination of the above loss functions:

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{self}} + \mathcal{L}_{\text{intra-inter}} + \mathcal{L}_{\text{sparse}} + \mathcal{L}_{\text{others}}. \quad (7)$$

Our designed losses $\mathcal{L}_{\text{intra-inter}}$ and $\mathcal{L}_{\text{sparse}}$ are equally incorporated with the basic objective $\mathcal{L}_{\text{self}}$. An overview of the training process is shown in Algorithm 1.

3.5 Online Inference

After offline training, the pretrained backbone model is fixed during the inference stage, which is utilized to compute the affinity matrix for label transformation (*e.g.*, segmentation mask). Note that the contrastive transformation is merely utilized for offline training, and the inference process is similar to the intra-video transformation. To acquire more reliable correspondence, we further design a mutually correlated affinity to exclude noisy matching as follows:

$$\tilde{\mathbf{A}}_{r \rightarrow t}(i, j) = \frac{\exp(\mathbf{w}(i, j) \mathbf{f}_t(i)^\top \mathbf{f}_r(j))}{\sum_j \exp(\mathbf{w}(i, j) \mathbf{f}_t(i)^\top \mathbf{f}_r(j))}, \quad (8)$$

where $\mathbf{w}(i, j) \in [0, 1]$ is a mutual correlation weight between two frames. Ideally, we prefer the one-to-one matching, *i.e.*, one pixel in the reference frame should be highly correlated with some pixel in the target frame and vice versa. The mutual correlation weight is formulated by:

$$\mathbf{w}(i, j) = \frac{\mathbf{f}_t(i)^\top \mathbf{f}_r(j)}{\max_{i \in [1, N_1]} (\mathbf{f}_t(i)^\top \mathbf{f}_r(j))} \times \frac{\mathbf{f}_t(i)^\top \mathbf{f}_r(j)}{\max_{j \in [1, N_2]} (\mathbf{f}_t(i)^\top \mathbf{f}_r(j))}. \quad (9)$$

The weight \mathbf{w} can be regarded as the affinity normalization across both reference and target spatial dimensions. Given the above affinity between two frames, the target frame label $\hat{\mathbf{L}}_t$ can be transformed via $\hat{\mathbf{L}}_t = \tilde{\mathbf{A}}_{r \rightarrow t} \mathbf{L}_r$.

4 Experiments

We verify the effectiveness of our method on a variety of vision tasks including video object segmentation, visual object tracking, pose keypoint tracking, and human parts segmentation propagation¹.

¹The source code and pretrained model will be available at <https://github.com/594422814/ContrastCorr>

Algorithm 1: Offline Training Process

Input: Unlabeled video sequences.

Output: Trained weights for the backbone network.

```

1 for each mini-batch do
2   Extract deep features of the video frames;
3   Patch-level tracking to obtain matched feature pairs;
4   for each video in the mini-batch do
5     // Intra- and Inter-video transformations
6     Compute intra-video affinity  $\mathbf{A}_{r \rightarrow t}$  (Eq. 1);
7     Compute inter-video affinity  $\mathbf{A}_{r, \Sigma \rightarrow t}$  (Eq. 3);
8     Conduct intra- and inter-video transformations;
9     // Loss Computation
10    Compute intra-video self-supervision  $\mathcal{L}_{\text{self}}$ ;
11    Compute intra-inter consistency  $\mathcal{L}_{\text{intra-inter}}$ ;
12    Compute regularization terms  $\mathcal{L}_{\text{sparse}}$  and  $\mathcal{L}_{\text{others}}$ ;
13  end
14  Back-propagate all the losses in this mini-batch;
15 end
```

4.1 Experimental Details

Training Details. In our method, the patch-level tracking and frame transformations share a ResNet-18 backbone network (He et al. 2016) with the first 4 blocks for feature extraction. The training dataset is TrackingNet (Müller et al. 2018) with about 30k video. Note that previous works (Wang, Jabri, and Efros 2019; Li et al. 2019) use the Kinetics dataset (Zisserman et al. 2017), which is much larger in scale than TrackingNet. Our framework randomly crops and tracks the patches of 256×256 pixels (*i.e.*, patch-level tracking), and further yields a 32×32 intra-video affinity (*i.e.*, the network stride is 8). The batch size is 16. Therefore, each positive embedding contrasts with $15 \times (32 \times 32 \times 2) = 30720$ negative embeddings. Since our method considers pixel-level features, a small batch size also involves abundant contrastive samples. We first train the intra-video transformation (warm-up stage) for the first 100 epochs and then train the whole framework in an end-to-end manner for another 100 epochs. The learning rate of both two stages is 1×10^{-4} and will be reduced by half every 40 epochs. The training stage takes about one day on 4 Nvidia 1080Ti GPUs.

Inference Details. For a fair comparison, we use the same testing protocols as previous works (Wang, Jabri, and Efros 2019; Li et al. 2019) in all tasks.

4.2 Framework Effectiveness Study

In Table 1, we show ablative experiments of our method on the DAVIS-2017 validation dataset (Ponttuset et al. 2017). The evaluation metrics are Jaccard index \mathcal{J} and contour-based accuracy \mathcal{F} . As shown in Table 1, without the intra-video guidance, inter-video transformation alone for self-supervision yields unsatisfactory results due to overwhelming noisy/negative samples. With only intra-video transformation, our framework is similar to the previous approach (Li et al. 2019). By jointly employing both of these two transformations under an intra-inter consistency constraint, our method obtains obvious performance improvements of 3.2% in \mathcal{J} and 3.4% in \mathcal{F} . The sparsity term of inter-video affinity encourages the embedding separation and further

Intra-video Transformation	Inter-video Transformation	Sparsity Constraint	Mutual Correlation	\mathcal{J} (Mean)	\mathcal{F} (Mean)
$\mathcal{L}_{\text{self}} + \mathcal{L}_{\text{others}}$	$\mathcal{L}_{\text{intra-inter}}$	$\mathcal{L}_{\text{sparse}}$			
✓				55.8	60.3
✓	✓			59.0	63.7
✓	✓	✓		59.2	64.0
✓	✓	✓	✓	60.5	65.5

Table 1: Analysis of each component of our method on the DAVIS-2017 validation dataset.

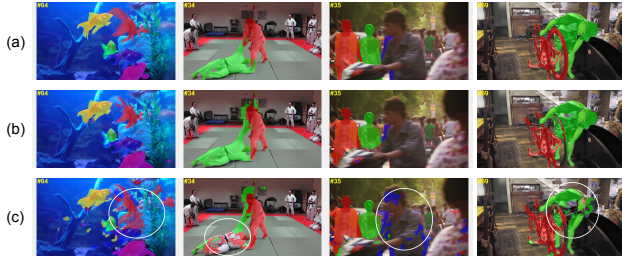


Figure 5: (a) Ground-truth results. (b) Results of the model with both intra- and inter-video transformations. (c) Results of the model without inter-video contrastive transformation, where the failures are highlighted by white circles.

improves the results.

In Figure 9, we further visualize the comparison results of our method with and without contrastive transformation. As shown in the last row of Figure 9, only intra-video self-supervision fails to effectively handle the challenging scenarios with distracting objects and partial occlusion. By involving the contrastive transformation, the learned feature embeddings exhibit superior discrimination capability for instance-level separation.

4.3 Comparison with State-of-the-art Methods

Video Object Segmentation on the DAVIS-2017. DAVIS (Ponttuset et al. 2017) is a video object segmentation (VOS) benchmark. We evaluate our method on the DAVIS-2017 validation set following Jaccard index \mathcal{J} (IoU) and contour-based accuracy \mathcal{F} . Table 2 lists quantitative results. Our model performs favorably against the state-of-the-art self-supervised methods including Time-Cycle (Wang, Jabri, and Efros 2019), CorrFlow (Lai and Xie 2019), and UVC (Li et al. 2019). Specifically, with the same experimental settings (e.g., frame input size and recurrent reference strategy), our model surpasses the recent top-performing UVC approach by 3.8% in \mathcal{J} and 4.8% in \mathcal{F} . The recent MAST approach (Lai, Lu, and Xie 2020) obtains impressive results by leveraging a memory mechanism, which can be added to our framework for further performance improvement. From Figure 8 (first row), we can observe that our method is robust in handling distracting objects and partial occlusion.

Compared with the fully-supervised ResNet-18 network trained on ImageNet with classification labels, our method exhibits much better performance. It is also worth noting that our method even surpasses the recent fully-supervised

Model	Supervised	\mathcal{J} (Mean)	\mathcal{F} (Mean)
Transitive Inv. (Wang, He, and Gupta 2017)		32.0	26.8
DeepCluster (Caron et al. 2018)		37.5	33.2
Video Colorization (Vondrick et al. 2018)		34.6	32.7
Time-Cycle (Wang, Jabri, and Efros 2019)		41.9	39.4
CorrFlow (Lai and Xie 2019)		48.4	52.2
UVC (480p) (Li et al. 2019)		56.3	59.2
UVC (560p) (Li et al. 2019)		56.7	60.7
MAST (Lai, Lu, and Xie 2020)		63.3	67.6
ContrastCorr (Ours)		60.5	65.5
ResNet-18 (He et al. 2016)	✓	49.4	55.1
OSVOS (Caelles et al. 2017)	✓	56.6	63.9
FEEVOS (Voigtlaender et al. 2019)	✓	69.1	74.0

Table 2: Evaluation on video object segmentation on the DAVIS-2017 validation dataset. The evaluation metrics are region similarity \mathcal{J} and contour-based accuracy \mathcal{F} .

Model	Supervised	DP@20pixel	AUC
KCF (HOG feature) (Henriques et al. 2015)		69.6	48.5
UL-DCFNet (Yang, Zhang, and Zhang 2019)		75.5	58.4
UDT (Wang et al. 2019)		76.0	59.4
UVC (Li et al. 2019)		-	59.2
LU DT (Wang et al. 2020)		76.9	60.2
ContrastCorr (Ours)		77.2	61.1
ResNet-18 + DCF (He et al. 2016)	✓	49.4	55.6
SiamFC (Bertinetto et al. 2016)	✓	77.1	58.2
DiMP-18 (Bhat et al. 2019)	✓	87.1	66.2

Table 3: Evaluation on video object tracking on the OTB-2015 dataset. The evaluation metrics are distance precision (DP) and area-under-curve (AUC) score of the success plot.

methods such as OSVOS.

Video Object Tracking on the OTB-2015. OTB-2015 (Wu, Lim, and Yang 2015) is a visual tracking benchmark with 100 challenging videos. We evaluate our method on OTB-2015 under distance precision (DP) and area-under-curve (AUC) metrics. Our model learns robust feature representations for fine-grained matching, which can be combined with the correlation filter (Henriques et al. 2015; Danelljan et al. 2014) for robust tracking. Without online fine-tuning, we integrate our model into a classic tracking framework based on the correlation filter, i.e., DCFNet (Wang et al. 2017). The comparison results are shown in Table 3. Note that UDT (Wang et al. 2019) is the recently proposed unsupervised tracker trained with the correlation filter in an end-to-end manner. Without end-to-end optimization, our model is still robust enough to achieve superior performance in comparison with UDT. Our method also outperforms the classic fully-supervised trackers such as SiamFC. As shown in Figure 8 (second row), our model can well handle the motion blur, deformation, and similar distractors.

Pose Keypoint Propagation on the J-HMDB. We evaluate our model on the pose keypoint propagation task on the validation set of J-HMDB (Jhuang et al. 2013). Pose keypoint tracking requires precise fine-grained matching, which is more challenging than the box-level or mask-level propagation in the VOT/VOS tasks. Given the initial frame with

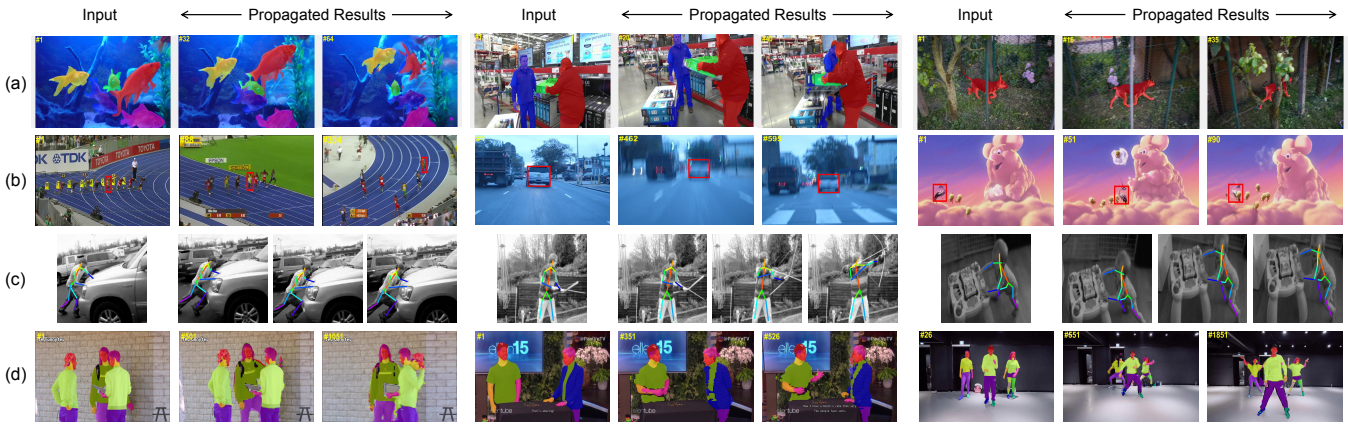


Figure 6: Experimental results of our method. (a) Video object segmentation on the DAVIS-2017. (b) Visual object tracking on the OTB-2015. (c) Pose keypoint tracking on the J-HMDB. (d) Parts segmentation propagation on the VIP.

Model	Supervised	PCK@.1	PCK@.2
SIFT Flow (Liu, Yuen, and Torralba 2011)		49.0	68.6
Transitive Inv. (Wang, He, and Gupta 2017)		43.9	67.0
DeepCluster (Caron et al. 2018)		43.2	66.9
Video Colorization (Vondrick et al. 2018)		45.2	69.6
Time-Cycle (Wang, Jabri, and Efros 2019)		57.3	78.1
CorrFlow (Lai and Xie 2019)		58.5	78.8
UVC (Li et al. 2019)		58.6	79.8
ContrastCorr (Ours)		61.1	80.8
ResNet-18 (He et al. 2016)	✓	53.8	74.6
Thin-Slicing Network (Song et al. 2017)	✓	68.7	92.1

Table 4: Keypoints propagation on J-HMDB. The evaluation metric is PCK at different thresholds.

15 annotated human keypoints, we propagate them in the successive frames. The evaluate metric is the probability of correct keypoint (PCK), which measures the percentage of keypoints close to the ground-truth in different thresholds. We show comparison results against the state-of-the-art methods in Table 4 and qualitative results in Figure 8 (third row). Our method outperforms all previous self-supervised methods such as Time-Cycle, CorrFlow, and UVC (Table 4). Furthermore, our approach significantly outperforms pre-trained ResNet-18 with ImageNet supervision.

Semantic and Instance Propagation on the VIP. Finally, we evaluate our method on the Video Instance-level Parsing (VIP) dataset (Zhou et al. 2018), which includes dense human parts segmentation masks on both the semantic and instance levels. We conduct two tasks in this benchmark: semantic propagation and human part propagation with instance identity. For the semantic mask propagation, we propagate the semantic segmentation maps of human parts (*e.g.*, heads, arms, and legs) and evaluate performance via the mean IoU metric. For the part instance propagation task, we propagate the instance-level segmentation of human parts (*e.g.*, different arms of different persons) and evaluate performance via the instance-level human parsing metric: mean Average Precision (AP). Table 5 shows that our method per-

Model	Supervised	mIoU	AP ^r _{vol}
SIFT Flow (Liu, Yuen, and Torralba 2011)		21.3	10.5
Transitive Inv. (Wang, He, and Gupta 2017)		19.4	5.0
DeepCluster (Caron et al. 2018)		21.8	8.1
Time-Cycle (Wang, Jabri, and Efros 2019)		28.9	15.6
UVC (Li et al. 2019)		34.1	17.7
ContrastCorr (Ours)		37.4	21.6
ResNet-18 (He et al. 2016)	✓	31.8	12.6
FGFA (Zhu et al. 2017)	✓	37.5	23.0
ATEN (Zhou et al. 2018)	✓	37.9	24.1

Table 5: Evaluation on propagating human part labels in Video Instance-level Parsing (VIP) dataset. The evaluation metrics are semantic propagation with mIoU and part instance propagation in AP^r_{vol}.

forms favorably against previous self-supervised methods. For example, our approach outperforms the previous best self-supervised method UVC by 3.3% mIoU in semantic propagation and 3.9% in human part propagation. Besides, our model notably surpasses the ResNet-18 model trained on ImageNet with classification labels. Finally, our method is comparable with the fully-supervised ATEN algorithm (Zhou et al. 2018) designed for this dataset.

5 Conclusion

In this work, we focus on the correspondence learning using unlabeled videos. Based on the well-studied intra-video self-supervision, we go one step further by introducing the inter-video transformation to achieve contrastive embedding learning. The proposed contrastive transformation encourages embedding discrimination while preserving the fine-grained matching characteristic among positive embeddings. Without task-specific fine-tuning, our unsupervised model shows satisfactory generalization on a variety of temporal correspondence tasks. Our approach consistently outperforms previous self-supervised methods and is even comparable with the recent fully-supervised algorithms.

Acknowledgements. The work of Wengang Zhou was supported in part by the National Natural Science Foundation of China under Contract 61822208, Contract U20A20183, and Contract 61632019; and in part by the Youth Innovation Promotion Association CAS under Grant 2018497. The work of Houqiang Li was supported by NSFC under Contract 61836011.

References

- Bertinetto, L.; Valmadre, J.; Henriques, J. F.; Vedaldi, A.; and Torr, P. H. 2016. Fully-convolutional siamese networks for object tracking. In *ECCV Workshop*.
- Bhat, G.; Danelljan, M.; Van Gool, L.; and Timofte, R. 2019. Learning discriminative model prediction for tracking. In *ICCV*.
- Caelles, S.; Maninis, K.; Ponttuset, J.; Lealtaix, L.; Cremers, D.; and Van Gool, L. 2017. One-shot video object segmentation. In *CVPR*.
- Caron, M.; Bojanowski, P.; Joulin, A.; and Douze, M. 2018. Deep clustering for unsupervised learning of visual features. In *ECCV*.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. E. 2020. A simple framework for contrastive learning of visual representations. *arXiv: 2002.05709*.
- Danelljan, M.; Häger, G.; Khan, F.; and Felsberg, M. 2014. Accurate scale estimation for robust visual tracking. In *BMVC*.
- Den Oord, A. V.; Li, Y.; and Vinyals, O. 2018. Representation learning with contrastive predictive coding. *arXiv: 1807.03748*.
- Dosovitskiy, A.; Fischery, P.; Ilg, E.; Hausser, P.; Hazirbas, C.; Golkov, V.; Der Smagt, P. V.; Cremers, D.; and Brox, T. 2015. FlowNet: learning optical flow with convolutional networks. In *ICCV*.
- Hadsell, R.; Chopra, S.; and Lecun, Y. 2006. Dimensionality reduction by learning an invariant mapping. In *CVPR*.
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In *CVPR*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*.
- Henriques, J. F.; Caseiro, R.; Martins, P.; and Batista, J. 2015. High-speed tracking with kernelized correlation filters. *TPAMI* 37(3): 583–596.
- Hjelm, R. D.; Fedorov, A.; Lavoie-marchildon, S.; Grewal, K.; Bachman, P.; Trischler, A.; and Bengio, Y. 2019. Learning deep representations by mutual information estimation and maximization. In *ICLR*.
- Huang, J.; Dong, Q.; Gong, S.; and Zhu, X. 2020. Unsupervised Deep Learning via Affinity Diffusion. In *AAAI*.
- Jabri, A.; Owens, A.; and Efros, A. 2020. Space-time correspondence as a contrastive random walk. In *NeurIPS*.
- Jhuang, H.; Gall, J.; Zuffi, S.; Schmid, C.; and Black, M. J. 2013. Towards understanding action recognition. In *ICCV*.
- Lai, Z.; Lu, E.; and Xie, W. 2020. MAST: A Memory-Augmented Self-Supervised Tracker. In *CVPR*.
- Lai, Z.; and Xie, W. 2019. Self-supervised learning for video correspondence flow. In *BMVC*.
- Lee, H.-Y.; Huang, J.-B.; Singh, M.; and Yang, M.-H. 2017. Unsupervised representation learning by sorting sequences. In *ICCV*.
- Li, X.; Liu, S.; De Mello, S.; Wang, X.; Kautz, J.; and Yang, M.-H. 2019. Joint-task self-supervised learning for temporal correspondence. In *NeurIPS*.
- Lin, T.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: common objects in context. In *ECCV*.
- Liu, C.; Yuen, J.; and Torralba, A. 2011. Sift flow: dense correspondence across scenes and its applications. *TPAMI* 33(5): 978–994.
- Liu, S.; Zhong, G.; De Mello, S.; Gu, J.; Jampani, V.; Yang, M.; and Kautz, J. 2018. Switchable temporal propagation network. In *ECCV*.
- Meister, S.; Hur, J.; and Roth, S. 2018. Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In *AAAI*.
- Müller, M.; Bibi, A.; Giancola, S.; Al-Subaihi, S.; and Ghanem, B. 2018. Trackingnet: a large-scale dataset and benchmark for object tracking in the wild. In *ECCV*.
- Pathak, D.; Girshick, R.; Dollár, P.; Darrell, T.; and Hariharan, B. 2017. Learning Features by Watching Objects Move. In *CVPR*.
- Pathak, D.; Krahenbuhl, P.; Donahue, J.; Darrell, T.; and Efros, A. A. 2016. Context encoders: feature learning by inpainting. In *CVPR*.
- Ponttuset, J.; Perazzi, F.; Caelles, S.; Arbelaez, P.; Sorkine-hornung, A.; and Van Gool, L. 2017. The 2017 davis challenge on video object segmentation. *arXiv: 1704.00675*.
- Song, J.; Wang, L.; Van Gool, L.; and Hilliges, O. 2017. Thin-slicing network: a deep structured model for pose estimation in videos. In *CVPR*.
- Tung, H. F.; Tung, H.; Yumer, E.; and Fragkiadaki, K. 2017. Self-supervised learning of motion capture. In *NeurIPS*.
- Voigtlaender, P.; Chai, Y.; Schroff, F.; Adam, H.; Leibe, B.; and Chen, L. 2019. Feelvos: fast end-to-end embedding learning for video object segmentation. In *CVPR*.
- Vondrick, C.; Shrivastava, A.; Fathi, A.; Guadarrama, S.; and Murphy, K. 2018. Tracking emerges by coloring videos. In *ECCV*.
- Wang, N.; Song, Y.; Ma, C.; Zhou, W.; Liu, W.; and Li, H. 2019. Unsupervised deep tracking. In *CVPR*.
- Wang, N.; Zhou, W.; Song, Y.; Ma, C.; Liu, W.; and Li, H. 2020. Unsupervised deep representation learning for real-time tracking. *IJCV* 1–19.
- Wang, Q.; Gao, J.; Xing, J.; Zhang, M.; and Hu, W. 2017. Dcfnet: discriminant correlation filters network for visual tracking. *arXiv:1704.04057*.

Wang, X.; and Gupta, A. 2015. Unsupervised learning of visual representations using videos. In *ICCV*.

Wang, X.; He, K.; and Gupta, A. 2017. Transitive invariance for self-supervised visual representation learning. In *ICCV*.

Wang, X.; Jabri, A.; and Efros, A. A. 2019. Learning correspondence from the cycle-consistency of time. In *CVPR*.

Wu, Y.; Lim, J.; and Yang, M.-H. 2015. Object tracking benchmark. *TPAMI* 37(9): 1834–1848.

Yang, L.; Zhang, D.; and Zhang, L. 2019. Learning a visual tracker from a single movie without annotation. In *AAAI*.

Ye, M.; Zhang, X.; Yuen, P. C.; and Chang, S. 2019. Unsupervised embedding learning via invariant and spreading instance feature. In *CVPR*.

Zhou, Q.; Liang, X.; Gong, K.; and Lin, L. 2018. Adaptive temporal encoding network for video instance-level human parsing. In *ACM MM*.

Zhu, X.; Wang, Y.; Dai, J.; Yuan, L.; and Wei, Y. 2017. Flow-guided feature aggregation for video object detection. In *CVPR*.

Zisserman, A.; Carreira, J.; Simonyan, K.; Kay, W.; Zhang, B.; Hillier, C.; Vijayanarasimhan, S.; Viola, F.; Green, T. F. G.; Back, T.; et al. 2017. The kinetics human action video dataset. *arXiv: 1705.06950*.

A Inference Details

In the inference stage, we leverage the computed affinity matrix to transform different types of inputs, *e.g.*, segmentation masks and pose keypoints. Similar to Time-Cycle and UVC, we adopt the same recurrent inference strategy to propagate the ground-truth result from the first frame, as well as the predicted results from the preceding L frames onto the target frame. We average all $L + 1$ predictions to obtain the final propagated map. Following previous works, L is set to 1 for the VIP dataset and 7 for all the rest benchmarks. For fair comparisons, following Time-Cycle and UVC, we also use the k-NN propagation schema and set $k = 5$ for all tasks. More details can be found in the source code.

B Transformation Results

In Figure 7, we exhibit some examples of our tracked image pairs. In our framework, we first randomly crop a reference patch in the reference frame and then conduct the patch-level tracking to form a pair of matched images. As shown in Figure 7, the image pairs have similar contents, which facilitate further intra- and inter-video transformations. Thanks to the patch-level tracking, our image pairs contain the real target appearance changes (*e.g.*, person view/pose changes), which differs from conventional contrastive methods based on the manually designed rules (*e.g.*, flip and rotation) to form image pairs.

In Figure 7, we also show the inter-video transformation results of our approach. The transformed images yield almost identical contents in comparison with the target patch, which affirms that our affinity matrix achieves reliable correspondence matching.

C Additional VOS Results

In Figure 8, we show more results of our approach on the DAVIS-2017 validation dataset. From Figure 8, we can observe that our method is able to accurately propagate the segmentation masks in challenging scenarios.

UVC algorithm represents the current state-of-the-art self-supervised correspondence approach based on the intra-video transformation paradigm. In contrast, our method further exploits the inter-video level transformation to reinforce instance-level embedding discrimination. In Figure 9, we further compare our approach with UVC. As shown in Figure 9, compared with UVC, our approach better handles the challenging scenarios such as occlusion, deformation, and similar distractors.

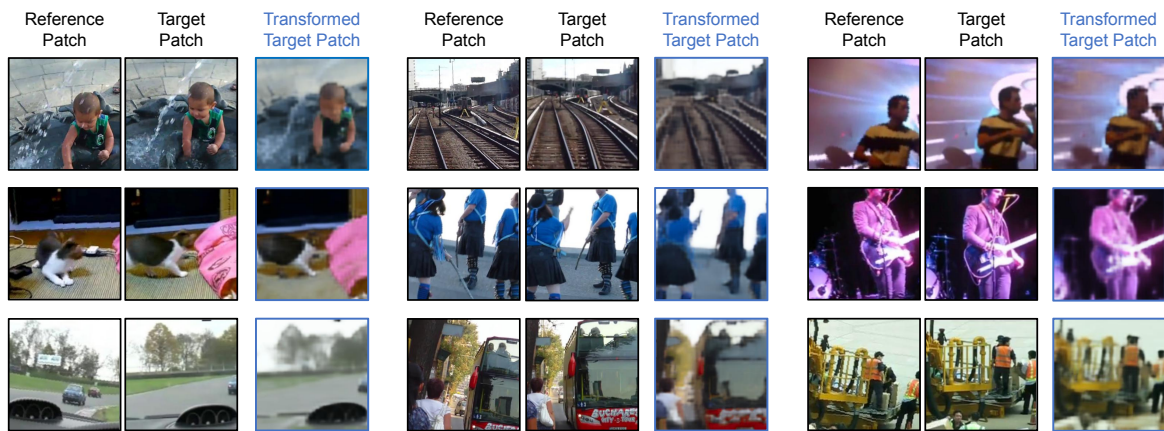


Figure 7: Examples of our tracked image pairs and transformed patches.

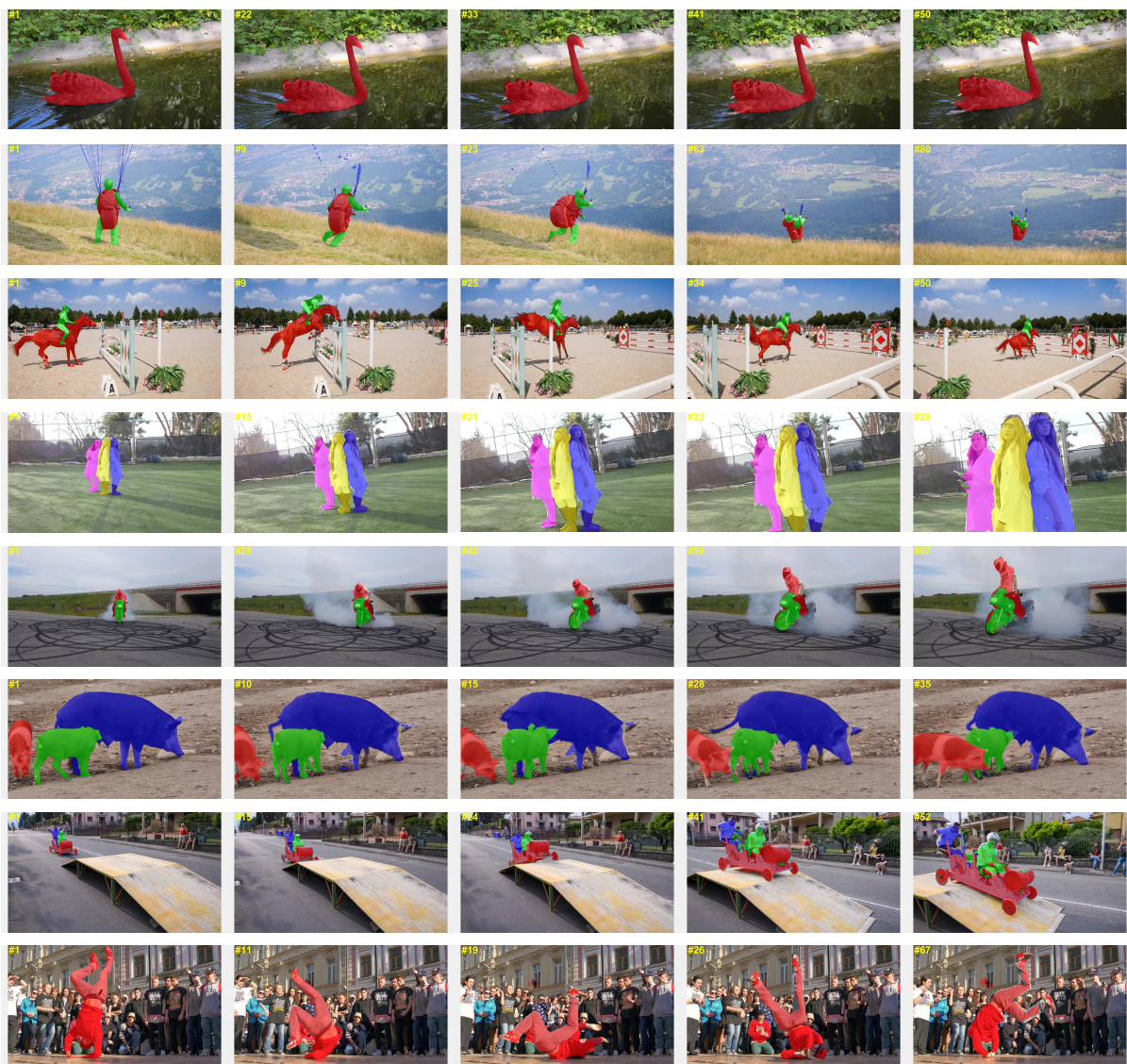


Figure 8: More results on the DAVIS-2017 validation dataset.

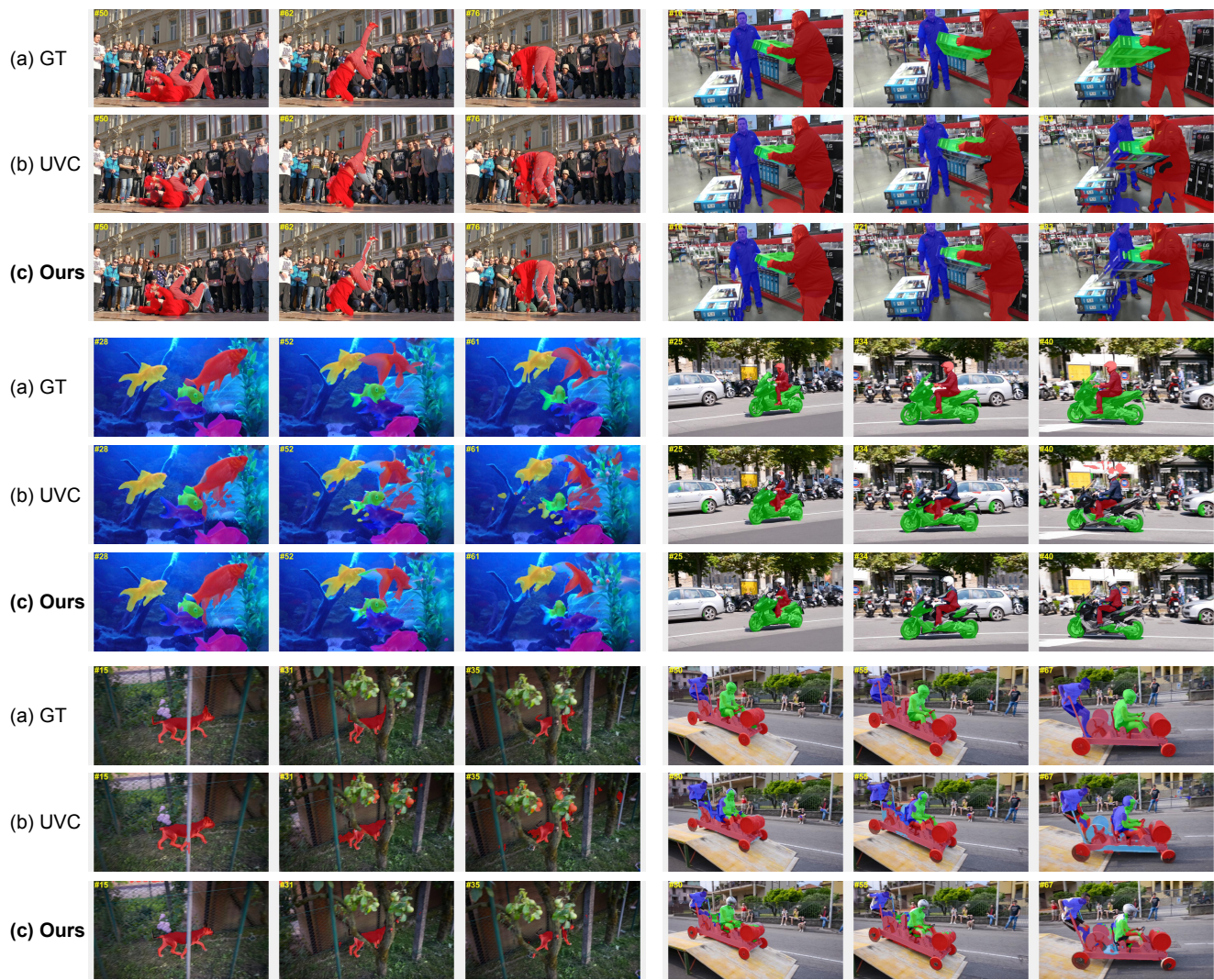


Figure 9: (a) Ground-truth segmentation results. (b) Results of UVC, which represents the current state-of-the-art performance of self-supervised correspondence methods. (c) Our results. By virtue of contrastive transformation, our approach shows superior results in comparison with previous intra-video based methods.