

# Self-supervised Learning for Video Correspondence Flow

Zihang Lai, Weidi Xie

## Contributions

Two main contributions of this work are:

- Simple information bottleneck (frame reconstruction by pixelwise matching) that forces the model to learn robust features for correspondence matching.
- The model is trained recursively on videos over long temporal windows to alleviate the tracker drifting effects with scheduled sampling and cycle consistency.

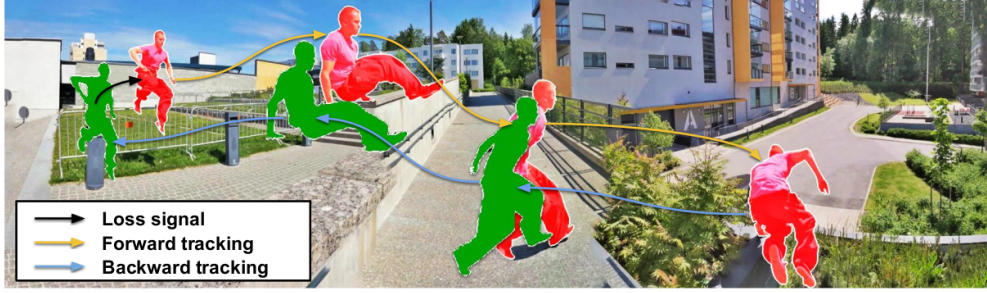


Figure 2: An overview of the proposed self-supervised learning for correspondence flow. A recursive model is used to compute the dense correspondence matching over a long temporal window with forward-backward cycle consistency.

## Method

**Feature Embedding:** Used ResNet as a feature encoder. Zero out 0, 1 or 2 channels in each RGB frame and apply data augmentation (brightness, contrast, saturation). Data augmentation prevents models from co-adaptation of low-level colors or illumination changes.

**Restricted Attention** It helps decrease in computation and memory consumption compared to full attention. Maximum disparity of  $M$  pixels is imposed in reference frame  $t$  to search for locally in square patch size  $(2M + 1) \times (2M + 1)$  centered at target pixel.

$$A^{ijkl} = \frac{\exp \left\langle f_t^{(i+k-M)(j+l-M)}, f_{t+1}^{ij} \right\rangle}{\sum_p \sum_q \exp \left\langle f_t^{(i+q)(j+p)}, f_{t+1}^{ij} \right\rangle} \quad (1)$$

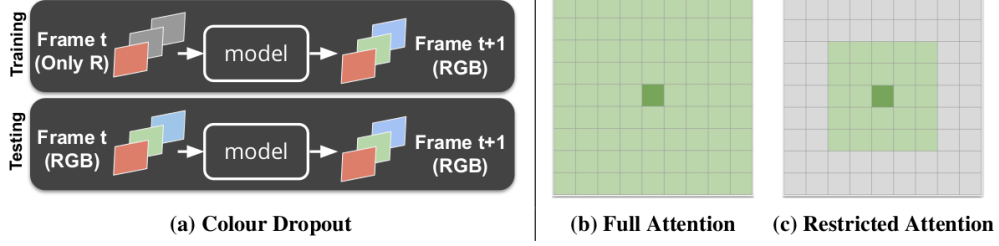


Figure 3: Restricted attention and colour dropout. See text for details.

where  $(i, j, k, l)$  is the entry of tensor denotes similarity between  $(i, j)$  of target frame, and pixel  $(i+k-M, j+l-M)$  of the reference frame

$$\hat{I}_{t+1} = \psi \left( A_{(t,t+1)}, I_t \right) = \sum_p \sum_q A^{ij(p+M)(q+M)} I_t \quad (2)$$

**Long-Term Correspondence Flow** Sampling training frames is difficult. Too close will result in no change and farther apart will result in way too much complex change.

- **Scheduled Sampling:** Replace Ground-truth tokens by model’s prediction. Shared embedding network is used to get feature embeddings ( $f_i = \Phi(g(I_i); \theta)$  where  $i = 1, \dots, n$ ). The reconstruction is a recursive process where nth frame ( $\tilde{I}_n$ ) may have access to previous frame’s ground truth ( $I_{n-1}$ ) or model prediction ( $\tilde{I}_{n-1}$ ).

$$\hat{I}_n = \begin{cases} \psi \left( A_{(n-1,n)}, I_{n-1} \right) \\ \psi \left( A_{(n-1,n)}, \hat{I}_{n-1} \right) \end{cases} \quad (3)$$

- **Cycle-Consistency:** It is used as a regularizer. Apply n frames forward and backward in future.

The final objective function is defined as:

$$L = \alpha_1 \cdot \sum_{i=1}^n \mathcal{L}_1 \left( I_i, \hat{I}_i \right) + \alpha_2 \cdot \sum_{j=n}^1 \mathcal{L}_2 \left( I_j, \hat{I}_j \right) \quad (4)$$

where  $\mathcal{L}_1, \mathcal{L}_2$  are pixelwise cross entropy loss between groundtruth and reconstructed frames in the forward and backward path.

## Results

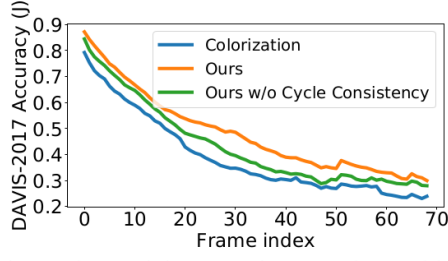


Figure 4: Model comparison on the problem of tracker drifting. The proposed model with cycle consistency has shown to be most robust as masks propagate.

Method	$\mathcal{J}(\text{Mean})$	$\mathcal{F}(\text{Mean})$
Ours (Full Model)	47.7	51.3
Ours w/o Colour Dropout	40.5	39.5
Ours w/o Restricted Attention	40.8	39.7
Ours w/o Scheduled Sampling	40.2	39.2
Ours w/o Cycle Consistency	41.0	40.4

Table 1: Ablation Studies on DAVIS-2017.  $\mathcal{J}$ : region overlapping,  $\mathcal{F}$ : contour accuracy respectively.

Method	Supervised	Dataset	$\mathcal{J} \& \mathcal{F}(\text{Mean})$	$\mathcal{J}(\text{Mean})$	$\mathcal{J}(\text{Recall})$	$\mathcal{F}(\text{Mean})$	$\mathcal{F}(\text{Recall})$
Identity	$\times$	-	22.9	22.1	15.9	23.6	11.7
Optical Flow (FlowNet2) [13]	$\times$	-	26.0	26.7	-	25.2	-
SIFT Flow [14]	$\times$	-	34.0	33.0	-	35.0	-
Transitive Inv. [15]	$\times$	-	29.4	32.0	-	26.8	-
DeepCluster [16]	$\times$	YFCC100M	35.4	37.5	-	33.2	-
Video Colorization [17]	$\times$	Kinetics	34.0	34.6	34.1	32.7	26.8
CycleTime (ResNet-50) [18]	$\times$	VLOG	40.7	41.9	40.9	39.4	33.6
<b>Ours (Full Model ResNet-18)</b>	$\times$	Kinetics [19]	<b>49.5</b>	<b>47.7</b>	<b>53.2</b>	<b>51.3</b>	<b>56.5</b>
<b>Ours (Full Model ResNet-18)</b>	$\times$	OxUvA [20]	<b>50.3</b>	<b>48.4</b>	<b>53.2</b>	<b>52.2</b>	<b>56.0</b>
ImageNet (ResNet-50) [19]	$\checkmark$	ImageNet	49.7	50.3	-	49.0	-
SiamMask [21]	$\checkmark$	YouTube-VOS	53.1	51.1	60.5	55.0	64.3
OSVOS[22]	$\checkmark$	DAVIS	60.3	56.6	63.8	63.9	73.8

Table 2: Video segmentation results on DAVIS-2017 dataset. Higher values are better.



(a) DAVIS 2017 Video Segmentation

(b) Keypoint Tracking

Method	Supervised	Dataset	$\text{PCK}_{\text{instance}}$		$\text{PCK}_{\text{max}}$	
			@.1	@.2	@.1	@.2
SIFT Flow[13]	$\times$	-	49.0	68.6	-	-
Video Colorization [17]	$\times$	Kinetics	45.2	69.6	-	-
CycleTime (ResNet-50) [18]	$\times$	VLOG	57.7	78.5	-	-
<b>Ours (Full Model ResNet-18)</b>	$\times$	Kinetics	<b>58.5</b>	<b>78.8</b>	<b>71.9</b>	<b>88.3</b>
ImageNet (ResNet-50) [19]	$\checkmark$	ImageNet	58.4	78.4	-	-
Fully Supervised [23]	$\checkmark$	JHMDB	-	-	68.7	81.6

Table 4: Keypoint tracking on JHMDB dataset (validation split 1). Higher values are better.