# Exploring the Semi-supervised Video Object Segmentation Problem from a Cyclic Perspective

**Yuxi Li · Ning Xu · Wenjie Yang · John See · Weiyao Lin**

**Abstract** Modern video object segmentation (VOS) algorithms have achieved remarkably high performance in a sequential processing order, while most of currently prevailing pipelines still show some obvious inadequacy like accumulative error, unknown robustness or lack of proper interpretation tools. In this paper, we place the semi-supervised video object segmentation problem into a cyclic workflow and find the defects above can be collectively addressed via the inherent cyclic property of semi-supervised VOS systems. Firstly, a cyclic mechanism incorporated to the standard sequential flow can produce more consistent representations for pixel-wise correspondance. Relying on the accurate reference mask in the starting frame, we show that the error propagation problem can be mitigated. Next, a simple gradient correction module, which naturally extends the offline cyclic pipeline to an online manner, can highlight the high-frequent and detailed part of results to further improve the segmentation quality while keeping feasi-

Yuxi Li
Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai China
E-mail: lyxok1@sjtu.edu.cn

Ning Xu
Adobe Research, San Jose, USA
E-mail: nxu@adobe.com

Wenjie Yang
Department of Computer Science, Shanghai Jiao Tong University, Shanghai China
E-mail: 13633491388@sjtu.edu.cn

John See
Heriot-Watt University, Malaysia
E-mail: J.See@hw.ac.uk

Weiyao Lin
Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai China
E-mail: wylin@sjtu.edu.cn

ble computation cost. Meanwhile such correction can protect the network from severe performance degration resulted from interference signals. Finally we develop cycle effective receptive field (cycle-ERF) based on gradient correction process to provide a new perspective into analyzing object-specific regions of interests. We conduct comprehensive comparison and detailed analysis on challenging benchmarks of DAVIS16, DAVIS17 and Youtube-VOS, demonstrating that the cyclic mechanism is helpful to enhance segmentation quality, improve the robustness of VOS systems, and further provide qualitative comparison and interpretation on how different VOS algorithms work. The code of this project can be found at `https://github.com/lyxok1/STM-Training`. [1]

## 1 Introduction

Video object segmentation (VOS) is garnering more attention in recent years due to its widespread application in the area of video editing and analysis. Among all the VOS scenarios, semi-supervised video object segmentation is the most practical and widely researched.

---

[1] This manuscript is an extended version of our conference paper to be published at the Thirty-fifth Conference on Neural Information Processing Systems (NeurIPS) 2020. *delving into the cyclic mechanism of semi-supervised video object segmentation* [15]. We have cited this paper in the manuscript and extended the paper substantially but not limited in following aspects: (1). A smooth regularized term and insight from frequency domain is appended in the method part. (2). In depth analysis on the robust of VOS model and the effect from our correction methods is included in the methods part. (3) More comprehensive experiments are appended to demonstrate the generality of our studies, including comparison under different baseline models and backbones, results with COCO pretraining and more qualitative results of effect from core components.

Specifically, a mask is provided in the first frame indicating the location and boundary of the objects, and the algorithm should accurately segment the same objects from the background in subsequent frames. A natural solution toward this problem is to process videos in a sequential order, exploiting the information from previous frames and guides the segmentation process in the current frame, since in most practical scenarios, the video is obtained in an online manner where only previous knowledge is available. Following this manner, current state-of-the-ar pipelines [34, 16, 18, 20, 31, 37, 27, 14] achieve high segmentation quality by delving into the information reuse from previous frames.

However, few research consider the flaws exposed by such sequential processing paradigm, where currently prevailing VOS pipelines still exhibit the following problems: (1) **Prone to accumulative error**. Ideally, if the masks predicted for intermediate frames are sufficiently accurate, they can provide more helpful object-specific information. Nevertheless, erroneous intermediate masks can mislead the segmentation procedure in future frames (as exemplified in Figure 1), and this error is further enlarged when low-quality segmentation dominate the reference templates. (2) **Robustness**. Although there is no research explicitly analyzing the robustness of VOS systems, but intuitively, it can be easily affected by noise manually appended into the previous knowledge, in cases such as adversarial attacks [7] particularly aimed at VOS algorithms. (3) **Unavailable tools for interpretation.** There is no unified tool to generally show how a VOS network is working given input frame and reference knowledge.

In this paper, we consider the problems above in a cyclical context and find the cyclic mechanism can be a potentially unified solution to collectively address these issues. Semi-supervised VOS is inherently suitable to be combined with a cyclic manner. Different from the predicted reference masks, the initial reference mask provided in the starting frame is always perfectly accurate and reliable. This inspires us to explicitly bridge the relationship between the initial reference mask and objective frame by taking the first reference mask as a measurement of prediction.

Specifically, when applying a generalized forward-backward data flow to form a cyclical structure and training our segmentation network at both the objective frame and starting frame, our model can learn more consistent correspondence relationship between predictions and the initial template mask. Further, at the inference stage, such cyclic structure can be naturally extended to an online version via a gradient correction module, which selectively refines the detailed and high-frequent part of predicted mask based on the gradient backward

from the starting frame at a marginal time cost. As a results, the cyclic workflow can effectively suppress the accumulative error with the starting frame as correct measurement. Meanwhile the online correction can also prevent the network from interference of noisy reference, boosting the robustness of our pipeline. Additionally, inspired by the process of gradient correction, we develop a new interpretation tool called cycle effective receptive field (cycle-ERF), which gradually updates an empty objective mask to show the strong response area w.r.t. the reference mask. In our experiments, we utilize the cycle-ERF to analyze how the cyclic training scheme affects the support regions of objects and highlight the difference in focus-area among distinctive baseline methods. This visualization method provides a fresh perspective for analyzing how the segmentation network extracts regions of interests from guidance masks.

The trained models are evaluated in both online and offline schemes on common object segmentation benchmarks: DAVIS16 [23], DAVIS17 [24] and Youtube-VOS [35], where we combine our cyclic methods to other baseline models and achieve results that are competitive to other state-of-the-art methods under a fair comparison setting. Besides, we also make detailed and comprehensive analysis to show how each part of the cyclic mechanism works under our design.

In a nutshell, the contributions of this paper can be summarized as follows:

- We incorporate cycle consistency into the training process of a semi-supervised video object segmentation network to mitigate the error propagation problem and further improve the segmentation quality. We achieved competitive results without data pretraining on mainstream benchmarks and can further be improved with more synthetic data.
- We design a gradient correction module to extend the offline segmentation network to an online approach, which boosts the model performance with marginal increase in computation cost, while keeping the model robust to disturbing noise.
- We develop cycle-ERF, a new visualization method to analyze the important regions on different segmentation models, which offers interpretability on the impact of cyclic training.

## 2 Related works

### 2.1 Semi-supervised video object segmentation

Semi-supervised video object segmentation has been widely researched in recent years with the rapid development of deep learning techniques. Depending on the

Fig. 1: An example of error propagation risk during the inference time, while the reference object is the camel in foreground, the distracting camel from background is incorrectly segmented at the same time.

presence of a learning process during inference stage, the segmentation algorithms can be generally divided into *online* methods and *offline* methods. OVOS [3] is the first online approach to exploit deep learning for the VOS problem, where a multi-stage training strategy is design to gradually shrink the focus of network from general objects to the one in reference masks. Subsequently, OnAVOS [32] improved the online learning process with an adaptive mechanism. MaskTrack [22] introduced extra static image data with mask annotation and employed data synthesized through affine transformation, to fine-tune the network before inference. All of these online methods require explicit parameter updating during inference. Although high performance can be achieved, these methods are usually time-consuming with a real-time FPS of less than 1, rendering them unfeasible for practical deployment.

On the other hand, there are a number of offline methods that are deliberately designed to learn generalized correspondence feature and they do not require necessary online learning process during inference time. RGMP [34] designed an hourglass structure with skip connections to predict the objective mask based on the current frame and previous information. S2S [35] proposed to model video object segmentation as a sequence-to-sequence problem and proceeds to exploit a temporal modeling module to enhance the temporal coherency of mask propagation. Other works like [37,18] resorted to using state-of-the-art instance segmentation or tracking pipeline [8,21] while attempting to design matching strategies to associate the mask over time. A few recent methods FEELVOS [31] and AGSS-VOS [16] mainly exploited the guidance from the initial reference and the last previous frame to enhance the segmentation accuracy with deliberately designed feature matching scheme or attention mechanism. STM [20] further optimized the feature matching process with external feature memory and an attention-based matching strategy, such memorial structure is further optimized by involving global context [14], adaptive gaussian kernel [27] and structure

of vision transformer [36]. Compared with online methods, these offline approaches are more efficient. However, to learn more general and robust pixel-wise feature correspondence, these data-hungry methods may require backbones pretrained on large amounts of extra data with mask annotations from other tasks such as instance segmentation [17,6] or saliency detection [28]. Without these auxiliary help, the methods might well be disrupted by distractions from similar objects in the video, which then propagates erroneous mask information to future frames.

All the approaches mentioned above follow a standard sequential processing order from start to end of the video and can not ensure the predicted mask is closely related to the initial reference guidance at first frame. In contrast, by embedding the cyclic mechanism into training stage, our method explicitly impose the constraint of reference mask in learning process. Besides, the online extension of gradient correction does not update the pretrained model parameters as other online learning methods, but dynamically refine the output according to the modification information from the reliable initial reference mask.

## 2.2 Cycle consistency

Cycle consistency is widely researched in unsupervised and semi-supervised representation learning, where a transformation and its inverse operation are applied sequentially on input data, the consistency requires that the output representation should be close to the original input data in feature space. With this property, cycle consistency can be applied to different types of correspondence-related tasks. [4] is a classical technique for correspondence learning, which treat the learning problem as video colorization and impose the consistency on natural color space to force embeddings of the same semantics to be closer in feature space. In the work of [29], a multiple step cycle-loop is build to con-

nect correspondence relationship between real images and rendered shots from 3D models. [33] combined patch-wise consistency with a weak tracker to construct a forward-backward data loop and this guides the network to learn representative feature across different time spans, [10] further extend such cyclic data loop to a random walk process. [19] exploited the cycle consistency in unsupervised optical flow estimation by designing a bidirectional consensus loss during training. On the other hand Cycle-GAN [39] and Recycle-GAN [2] and other popular examples of how cyclic training can be utilized to learn non-trivial cross-domain mapping, yielding reliable image-to-image transformation across different domains.

Our method with cyclic mechanism is different from the works mentioned above in following aspects. First, the motivation to exploit cycle consistency in our work is to explicitly regularize the predicted mask to be accurate for backward reference in cyclic loop, while in unsupervised methods like [33], the cycle consistency is an approach to obtain correspondence ground-truth. Further, the learning objective of our work is still a fully supervised segmentation task, with high-level and clear semantic information (the object is taken as foreground while the others are background). In contrast, methods with unsupervised cycle consistency for correspondence learning construct their objective by self-mimic in low-level semantics (e.g. the color space [4], spatial position [33] or transformation flow [29]), thus the learned embeddings are easy to correspond to distractors with similar low-level expression if there is not sufficient context provided. Consequently large amount of data are required to train these unsupervised methods to learn to catch the context relationship. Finally, our cyclic structure is not only applicable during training, but also useful in the inference stage. By measuring the consistency between initial reference mask and predicted results, we can refine the output on current frame to obtain more accurate guidance for future prediction.

## 3 Methods

### 3.1 Problem formulation

Given a video of length $T$, $X_t$ is the $t$-th frame ($t \in [1, T]$) in temporal sequential order, and $Y_t$ is its corresponding annotation mask. $\mathcal{S}_\theta$ is an object segmentation network parameterized by learnable weights $\theta$. In terms of the sequential processing order of the video, the segmentation network should achieve the function as in Equation (1) below:

$$\widehat{Y}_t = \mathcal{S}_\theta \left( \mathcal{X}_{t-1}, \mathcal{Y}_{t-1}, X_t \right) \quad t \in [2, T] \tag{1}$$

where $\widehat{Y}_t$ denotes the predicted object mask at $t$-th frame. $\mathcal{X}_{t-1} \subset \{X_i | i \in [1, t-1]\}$ is the reference frame set, which is a subset of all frames appearing before objective frame $X_t$. Similarly, $\mathcal{Y}_{t-1}$ is a set containing reference object masks corresponding to the reference frames in $\mathcal{X}_{t-1}$. However, in the semi-supervised setting, only the initial reference mask at the first frame is available. Therefore, in the inference stage, the corresponding predicted mask $\widehat{Y}_t$ is taken as the approximation of the reference mask. Hence, we have $\mathcal{Y}_{t-1} \subset \{Y_1\} \bigcup \{\widehat{Y}_i | i \in [2, t-1]\}$.

### 3.2 Cycle consistency loss

For the sake of mitigating error propagation during training, we incorporate the cyclical process into the offline training process to explicitly bridge the relationship between the initial reference and predicted masks. To be specific, as illustrated in Figure 2, after obtaining the predicted output mask $\widehat{Y}_t$ at frame $t$, we construct a **cyclic reference set** for frames and mask set, respectively denoted as $\widehat{\mathcal{X}}_t \subset \{X_i | i \in [2, t]\}$, $\widehat{\mathcal{Y}}_t \subset \{\widehat{Y}_i | i \in [2, t]\}$.

With the cyclic reference set, we can obtain the prediction for the initial reference mask in the same manner as sequential processing:

$$\widehat{Y}_1 = \mathcal{S}_\theta \left( \widehat{\mathcal{X}}_t, \widehat{\mathcal{Y}}_t, X_1 \right) \tag{2}$$

Consequently, we apply mask reconstruction loss (in Equation 3) during supervision, optimizing on both the output mask of the $t$-th frame $\widehat{Y}_t$ and the backward prediction $\widehat{Y}_1$.

$$\mathcal{L}_{cycle,t} = \mathcal{L}(\widehat{Y}_t, Y_t) + \mathcal{L}(\widehat{Y}_1, Y_1) \tag{3}$$

In implementation, we utilize the combination of cross-entropy loss and mask IOU loss as supervision at both sides of the cyclic loop, which can be formulated as,

$$\mathcal{L}(\widehat{Y}_t, Y_t) = \mathcal{L}_{IOU} + \gamma \mathcal{L}_{CE} \tag{4}$$

$$\mathcal{L}_{IOU} = 1 - \frac{\sum_{u \in \Omega} \min(\widehat{Y}_{t,u}, Y_{t,u})}{\sum_{u \in \Omega} \max(\widehat{Y}_{t,u}, Y_{t,u})} \tag{5}$$

$$\mathcal{L}_{CE} = \sum_{u \in \Omega} \left( (1 - Y_{t,u}) \log(1 - \widehat{Y}_{t,u}) + Y_{t,u} \log(\widehat{Y}_{t,u}) \right) \tag{6}$$

where $\Omega$ denotes the set of all pixel coordinates in the mask while $Y_{t,u}$ and $\widehat{Y}_{t,u}$ are the normalized pixel values at coordinate $u$ of the masks, $\gamma$ is a hyperparameter to balance between the two loss terms. It should also be noted that the cyclic mechanism in Figure 2 indirectly applies data augmentation on the training data
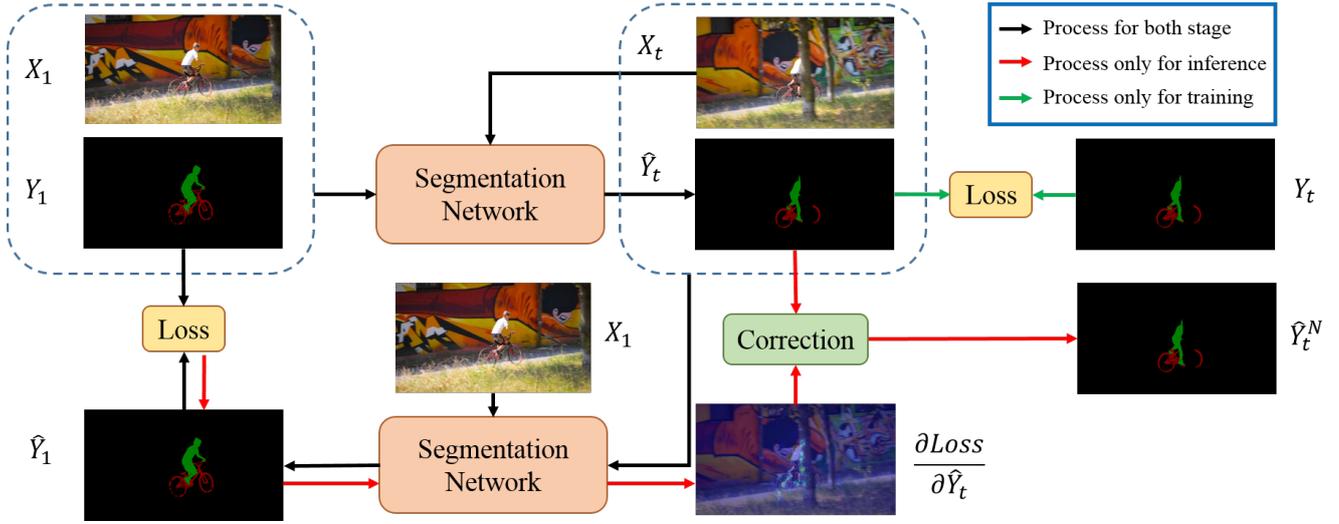
Fig. 2: Overview of the proposed cyclic mechanism in both training and inference stages of the segmentation network. For simplicity, we take the situation where $\mathcal{X}_{t-1} = \{X_1\}$, $\mathcal{Y}_{t-1} = \{Y_1\}$, $\widehat{\mathcal{X}}_t = \{X_t\}$ and $\widehat{\mathcal{Y}}_t = \{\widehat{Y}_t\}$ as an example.

by reversing the input clips in temporal order, helping the segmentation network to learn more general feature correspondences.

### 3.3 Gradient correction

**Cyclic refinement process.** After training with the cyclic loss as Equation (3), we can directly apply the offline model in the inference stage. However, inspired by the cyclic structure in training process, we can take the accurate initial reference mask as a measurement to evaluate the segmentation quality of current frame and proceed to refine the output results based on the evaluation results. In this way, we can explicitly reduce the effect of error propagation during inference time to keep our trained model from disturbances by natural or adversarial noises.

To achieve this goal, we design a gradient correction block to update segmentation results iteratively as illustrated in Figure 2. Since only the initial mask $Y_1$ is available in inference stage, we apply the predicted mask $\widehat{Y}_t$ to infer the initial reference mask in the same manner as Equation (2), and then evaluate the segmentation quality of $\widehat{Y}_t$ with the loss function in Equation (4). Intuitively, when more accurate prediction mask $\widehat{Y}_t$ are taken as reference, smaller reconstruction error for $Y_1$ will be yielded; therefore, during the gradient correction, our algorithm is focusing on minimizing the reconstruction error of the reference mask

$$\min_{\widehat{Y}_t} \mathcal{L}_{rec} = \min_{\widehat{Y}_t} \mathcal{L}\left(\mathcal{S}_\theta\left(\{X_t\}, \{\widehat{Y}_t\}, X_1\right), Y_1\right) \quad (7)$$

Where the loss term $\mathcal{L}(\cdot, \cdot)$ adopts the same formulation as Equation (4) The gradient descent method is adopted to solve the reconstruction problem in Equation (7) so as to refine the mask $\widehat{Y}_t$. To be specific, we start from an output mask $\widehat{Y}_t^0 = \widehat{Y}_t$, and then update the mask for $N$ iterations:

$$\widehat{Y}_t^{l+1} = \widehat{Y}_t^l - \alpha \frac{\partial \mathcal{L}_{rec}}{\partial \widehat{Y}_t^l} \quad (8)$$

where $\alpha$ is a predefined correction rate for mask update and $N$ is the iteration times. With this iterative refinement, we naturally extend the offline model to an online inference algorithm. However, the gradient correction approach can be time-consuming since it requires multiple times of network forward-backward pass. Due to this reason, we only apply gradient correction once per $K$ frames to achieve good performance-runtime trade-off.

**Interpretation in the frequency domain.** Empirically, with the gradient correction process in Equation (8), the details of output mask can be better handled. This claim can be demonstrated from the aspect of frequency domain, where we find that the gradient correction module empirically acts a high-frequency amplifier to polish the output mask in fine-grained details. To analyze from the aspect of frequency, we take the gradient correction module as a black box system and calculate its frequency response by computing the averaged ratio between output and input amplitude-frequency characteristics,

$$\mathcal{A}\mathcal{F}_{GC} = \frac{1}{T} \sum_{t=1}^{T} \frac{\left|\mathcal{F}\mathcal{F}\mathcal{T}\left(\widehat{Y}_t^N\right)\right|}{\left|\mathcal{F}\mathcal{F}\mathcal{T}\left(\widehat{Y}_t^0\right)\right|} \quad (9)$$
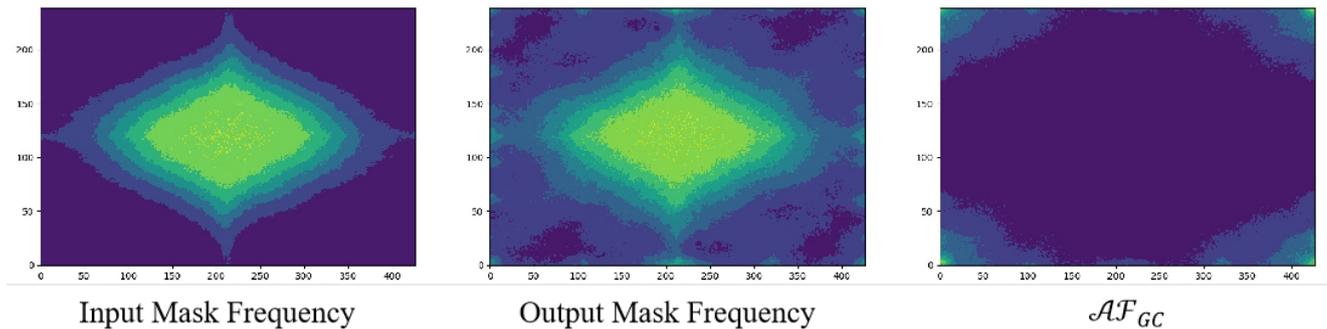
Fig. 3: Frequency domain distribution of input mask, output mask and amplitude response of gradient correction module. The results is obtained by averaging the response value on all frames of DAVIS17 validation set.

where $\mathcal{FFT}(\cdot)$ denotes 2D Fast Fourier transform. We visualized the 2D frequency-domain intensity of the input mask, output mask of gradient correction and the frequency-domain response $\mathcal{AF}_{GC}$ in the form of amplitude-frequency characteristic contour in Figure 3. We observe that compared with the input mask, the output mask manifests stronger intensity in high-frequency component (the four corner of the contour figure), the frequency response of gradient correction module also highlights and amplifies the part corresponding high-frequency harmonic. This observation provides the explanation of higher boundary accuracy improvement from the perspective of frequency since the high-frequency component usually supplements some detailed information of output masks.

**Anti-noise regularization.** However, intuitively, amplifying the high-frequency component will not only append details but also results in noise and artifacts around the object boundaries. To overcome such artifacts from gradient-correction, we augment the reconstruction error with a regularization term to suppress potential noise after refinement, denoted as

$$L_{rec} = \mathcal{L}\left(\mathcal{S}_\theta\left(\{X_t\}, \{\widehat{Y}_t\}, X_1\right), Y_1\right) + \lambda \mathcal{L}_{smooth}\left(\widehat{Y}_t\right) \tag{10}$$

where $\mathcal{L}_{smooth}\left(\widehat{Y}_t\right)$ is a spatial smooth term to avoid exaggerating some spot-like area in output masks with $\lambda$ as its corresponding weight parameter.

$$\mathcal{L}_{smooth}\left(\widehat{Y}_t\right) = \frac{1}{|\Omega|} \sum_{u \in \Omega} \nabla_x \widehat{Y}_{t,u}^2 + \nabla_y \widehat{Y}_{t,u}^2 \tag{11}$$

In instantiation, the Sobel operators $\nabla_x, \nabla_y$ are first applied on the mask $\widehat{Y}_t$ along the horizontal and vertical directions to calculate the spatial gradient, then areas with large spatial gradient norm are penalized. The augmented reconstruction objective can still be optimized according to the manner of Equation (8).

## 3.4 Cycle-ERF

The cyclic mechanism with gradient update in Equation (8) is not only helpful for the output mask refinement, but it also offers a new aspect of analyzing the region of interests of specific objects segmented by the pretrained network. In detail, we construct a reference set, $\mathcal{X}_l = \{X_l\}$ and $\mathcal{Y}_l = \{\mathbf{0}\}$ as the guidance, where $\mathbf{0}$ denotes an empty mask of the same size as $X_l$ but is filled with zeros. We take these references to predict objects at the $t$-th frame $\widehat{Y}_t$. To this end, we can obtain the prediction loss $\mathcal{L}(\widehat{Y}_t, Y_t)$. To minimize this loss, we conduct the gradient correction process as in Equation (8) to gradually update the empty mask for $M$ iterations. Finally, we take the ReLU function to preserve the positively activated areas of the objective mask as our final cycle-ERF representation, the resulting receptive field can be expressed as

$$\textbf{cycle-ERF}(Y_l) = ReLU\left(\widehat{Y}_l^M|_{\widehat{Y}_l^0 = \mathbf{0}}\right) \tag{12}$$

As we will show in our experiments, the cycle-ERF is capable of properly reflecting the support region of specific objects for the segmentation task. Through this analysis, the pretrained model can be shown to be particularly concentrated on certain objects in videos, making the learned segmentation model more interpretable.

## 4 Experiments

### 4.1 Experiment setup

**Datasets.** We train and evaluate our method on three widely used benchmarks for semi-supervised video object segmentation, DAVIS16 [23], DAVIS17 [24], and Youtube-VOS [35]. DAVIS16 contains 50 videos in total, where 30 sequences are used for training and the others are taken as validation. In this benchmark, each video

only covers a single reference mask. DAVIS17 is a multi-object extension of DAVIS16 and contains 120 video sequences in total with at most 10 objects in a video. The dataset is split into 60 sequences for training, 30 for validation, and the other 30 for test. The Youtube-VOS is larger in scale and contains more object categories. There are a total of 3,471 video sequences for training and 474 videos for validation in this dataset with at most 12 objects in a video, it also contains videos where objects appear from intermediate frames. Following the training procedure in [20,31], we construct a hybrid training set by mixing the data from all training sequences.

**Metrics.** For evaluation on DAVIS16, DAVIS17 validation, and test set, we adopt the metric following standard DAVIS evaluation protocol [24]. The Jaccard overlap $\mathcal{J}$ is adopted to evaluate the mean IOU between predicted and groundtruth masks. The contour F-score $\mathcal{F}$ computes the F-measurement in terms of the contour based precision and recall rate. The final score is obtained from the average value of $\mathcal{J}$ and $\mathcal{F}$. The evaluation on Youtube-VOS follows the same protocol except that the two metrics are computed on seen and unseen objects respectively and averaged together.

**Baselines.** We take three recent methods as our base model for further analysis of our proposed cyclic mechanism in the following experiments.

- Space Time Memory Network (STM) [20] is a widely used pipeline for fast and accurate semi-supervised video object segmentation, and the memory mechanism makes it flexible in adjusting reference sets $\mathcal{X}_t$ and $\mathcal{Y}_t$, which is suitable as comparison to our reference-based approach.
- Kernelized Memory Network (KMN) [27] is the kernelized version of STM network, where a Gaussian kernel is dynamically calculated before merging the knowledge from memory into current query feature.
- Global Context Memory Network (GCM) [14] extends the STM-like structure with a global temporal span, where the spatially global context of each reference frame is stored and dynamically updated in the memory.
- Associate Objects with Transformer (AOT) [36] is a more advanced and efficient video object segmentation framework, where a long-short term cross attention structure is designed to help with parallel video object segmentation on multiple objects.

Since there is no public training code of [20,27,14], we implement them by ourselves, and for [36], we directly implement our cycle mechanism from the public code from the author. For STM, in order to adapt to the time-consuming gradient correction process, we take the lightweight design by reducing the intermediate memory

feature dimension, resizing the input resolution for inference to $240 \times 427$, which is $1/4$ of the size in original work [20] ($480 \times 854$), and then upsampling the output to original resolution by nearest interpolation. For KMN, we replace the argmax operation in original paper with soft-argmax to make sure the network is end-to-end trainable. For AOT, we modify the one-hot encode into soft labels to ensure the cycle loop is end-to-end differentiable. For ease of representation, we denote the models trained with cyclic scheme with a "-cycle" suffix. It should be mentioned that the adopted baseline models from the original papers involved different external static data from various segmentation datasets [17,6], resulting in unfair and inconsistent comparisons here. To enable fairer and more consistent comparison, for most of our analysis, we re-train our implemented models using only the training data in DAVIS and Youtube-VOS. Nevertheless, for more comprehensive comparison, we still provide results of STM and AOT with the same setting as [20], where the model is pretrained purely on COCO [17] and predict the object mask at the resolution of ($480 \times 854$), which is the most commonly adopted experimental setting.

**Implementation details.** The training and inference procedures are deployed on an NVIDIA TITAN Xp GPU. Within an epoch, for each video sequence, we randomly sample 3 frames as the training samples – the frame with the smallest timestamp is regarded as the initial reference frame. Similar to [20], the maximum temporal interval of sampling increases by 5 every 20 training epochs. We set the hyperparameters as $\gamma = 1.0$, $\lambda = 0.75$, $N = 10$, $K = 5$, and $M = 50$. During training, we adopt a bootstrapping strategy for the cross entropy loss, where only the top $40\%$ pixels with maximum training loss are taken into account. The ResNet series of models [9] pretrained on ImageNet [26] are adopted as our backbone for baseline. The network is trained with a batch size of 4 for 240 epochs in total and is optimized by the Adam optimizer [13] of learning rate $10^{-5}$ and $\beta_1 = 0.9, \beta_2 = 0.999$. In both training and inference stages, the input frames are resized to the resolution of $240 \times 427$. The final output is upsampled to the original resolution by nearest interpolation. For simplicity, we directly use $X_t$ and $\widehat{Y}_t$ to construct the cyclic reference sets. For the case of multiple objects, we adopt the soft aggregation method in [16,20] to normalize the probability at each pixel of output masks.

**Data augmentation.** We apply common augmentation operations including random horizontal flipping, additive Gaussian noise and contrast enhancement. Additionally, we also adopt random crop strategy and fixed affine transformation (sheer, resize, rotation) to each training sample. We note that frames selected from the

same video will share the same transformation parameters. When exploiting COCO to synthesize data, we follow [1] to take random affine and copy-paste trick to generate pseudo sequences.

## 4.2 Main results

In this section, we first report the comparison results between our cyclic model and other methods, where our full model is measured with configuration of different backbone and cyclic schemes.

**DAVIS.** The evaluation results on DAVIS16 validation set, DAVIS17 validation and test-dev set are reported in Table 1. From this table, we observe that our model trained with cyclic loss outperforms most of the offline methods and even performs better than the method with online learning [32] on DAVIS17 benchmark. When combined with the online gradient correction process, our method gets further improvement. It should also be noticeable that although standard STM-cycle do not require additional training data other than DAVIS and Youtube-VOS, it outperforms some other state-of-the-art pipelines highly dependent on additional training data [16,34,31]. In terms of the runtime speed, although gradient correction increases the computation cost, our method still runs at a speed comparable to other offline methods [16] due to our efficient implementation. When we replace the backbone network from ResNet50 to more lightweight ResNet18, our trained model can run faster with still competitive performance. Although there is a performance gap between our approach that is trained from scratch and the state-of-the-art online learning method, our method is far more efficient and it does not requires collecting extra data from instance segmentation tasks as training samples. It is also noticeable that when adding COCO into training set, cyclic version of STM can achieve state-of-the-art performance on all benchmarks of DAVIS, obtaining consistent improvement over original STM [20], which shares the same backbone but requires more data besides COCO.

Furthermore, we also try to combine the gradient correction process with an existing open source model of AGSS-VOS [18][2] to test inference performance gain on state-of-the-art method with purely gradient correction. This appears to bring improvement on the overall segmentation quality, demonstrating that cyclic consistency is helpful even in different segmentation pipelines.

**Youtube-VOS.** The evaluation results on Youtube-VOS validation set are reported in Table 2. On this benchmark, our model also outperforms some offline methods and their online learning counterparts [35,37].

---

[2] `https://github.com/Jia-Research-Lab/AGSS-VOS`

It is also noticeable that compared with the performance on seen objects, the one on unseen objects has improved more using our gradient correction strategy. Further, we observe that with ResNet-50 backbone and gradient correction technique, our final pipeline can achieve similar segmentation accuracy and runtime speed to some state-of-the-art methods [16] on Youtube-VOS even without extra data for pre-training. Finally, as consistent with the results on DAVIS series, when combined with COCO pretraining, our STM-cycle model can achieves better performance than other state-of-the-art models with pretrained [20,16,34]

## 4.3 Ablation study

In this section, we conduct a series of experiments to analyze the cyclic property of our trained models, with all the results evaluated on the DAVIS16 and DAVIS17 validation set and ResNet50 as backbone network.

### 4.3.1 Effectiveness of each cyclic component.

We first demonstrate the effectiveness of cyclic training and gradient correction in Table 4, where the baseline method [20] based on STM network is re-implemented and retrained.From this table, both components are shown to be helpful in boosting the performance on both DAVIS16 and DAVIS17 validation sets. In particular, the incorporated cycle mechanism improves the contour score $\mathcal{F}$ more than the overlap score $\mathcal{J}$, signifying that the proposed scheme is likely to be more useful for fine-grained mask prediction.

### 4.3.2 Improvement with different reference sets.

Due to the flexibility of our baseline method in configuring its reference sets during inference, we tested how our cyclic training strategy would impact VOS performance using different reference sets on our STM baseline. We conduct the test under four types of configuration: (1) Only the initial reference mask and its frame are utilized for predicting other frames. (2) Only the prediction of the last frame $\widehat{Y}_{t-1}$ and the last frame are used. (3) Both the initial reference and last frame prediction are utilized, which is the most common configuration in other state-of-the-art works. (4) The external memory strategy (denoted as **MEM**) in [20] is used where the reference set is dynamically updated by appending new prediction and frames at a specific frequency of 5Hz. In the results reported in Table 3, we observe that the cyclic training is helpful under all configurations. It is also interesting to see that our scheme achieves the maximum improvement ($+4.6$ $\mathcal{J}\&\mathcal{F}$ on DAVIS17 and $+1.9$ $\mathcal{J}\&\mathcal{F}$ on DAVIS16)

| DAVIS17 validation | | | | | | | |
|---|---|---|---|---|---|---|---|
| Method | Extra data | OL | GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | FPS |
| RGMP [34] | ✓ | | | 64.8 | 68.6 | 66.7 | 3.6 |
| DMM-Net [37] | ✓ | | | 68.1 | 73.3 | 70.7 | - |
| AGSS-VOS [16] | ✓ | | | 63.4 | 69.8 | 66.6 | 10 |
| AGSS-VOS [16] | ✓ | | ✓ | 64.0 | 70.6 | 67.3 | 2.2 |
| FEELVOS [31] | ✓ | | | 69.1 | 74.0 | 71.5 | 2 |
| FRTM [25] | | | | - | - | 70.2 | 41.3 |
| OnAVOS [32] | ✓ | ✓ | | 61.0 | 66.1 | 63.6 | 0.04 |
| PReMVOS [18] | ✓ | ✓ | | 73.9 | 81.7 | 77.8 | 0.03 |
| STM-ResNet50 [20]$^{\dagger}$ | ✓ | | | 79.7 | 84.4 | 82.0 | 7.9 |
| STM-ResNet50-cycle (Ours)$^{\dagger}$ | ✓ | | ✓ | **80.4** | **85.2** | **82.8** | 2.5 |
| STM-ResNet18-cycle (Ours) | | | | 64.7 | 69.9 | 67.3 | **55.3** |
| STM-ResNet18-cycle (Ours) | | | ✓ | 65.3 | 70.8 | 68.1 | 13.4 |
| STM-ResNet50-cycle (Ours) | | | | 68.7 | 74.7 | 71.7 | 38 |
| STM-ResNet50-cycle (Ours) | | | ✓ | 69.8 | 75.9 | 72.9 | 9.3 |
| DAVIS17 test-dev | | | | | | | |
| Method | Extra data | OL | GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | FPS |
| RVOS [30] | | | | 48.0 | 52.6 | 50.3 | 22.7 |
| RGMP [34] | ✓ | | | 51.3 | 54.4 | 52.8 | 2.4 |
| AGSS-VOS [16] | ✓ | | | 54.8 | 59.7 | 57.2 | 10 |
| FEELVOS [31] | ✓ | | | 55.2 | 60.5 | 57.8 | 1.8 |
| OnAVOS [32] | ✓ | ✓ | | 53.4 | 59.6 | 56.9 | 0.03 |
| PReMVOS [18] | ✓ | ✓ | | 67.5 | 75.7 | 71.6 | 0.02 |
| STM-ResNet50 [20]$^{\dagger}$ | ✓ | | | 68.0 | 74.1 | 71.0 | 14.8 |
| STM-ResNet50-cycle (Ours)$^{\dagger}$ | ✓ | | ✓ | **70.6** | **76.4** | **73.5** | 4.1 |
| STM-ResNet18-cycle (Ours) | | | | 53.2 | 58.4 | 55.8 | **44.7** |
| STM-ResNet18-cycle (Ours) | | | ✓ | 53.7 | 60.5 | 57.2 | 10.7 |
| STM-ResNet50-cycle (Ours) | | | | 55.1 | 60.5 | 57.8 | 31 |
| STM-ResNet50-cycle (Ours) | | | ✓ | 55.4 | 62.8 | 59.1 | 6.9 |
| DAVIS16 validation | | | | | | | |
| Method | Extra data | OL | GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | FPS |
| Lucid Dreaming [12] | | | | 83.9 | 82.0 | 83.0 | - |
| RGMP [34] | ✓ | | | 81.5 | 82.0 | 81.8 | 7.8 |
| AGAME [11] | ✓ | | | 81.5 | 82.2 | 81.9 | 14.3 |
| FEELVOS [31] | ✓ | | | 81.1 | 82.2 | 81.7 | - |
| FRTM [25] | | | | - | - | 78.5 | 41.3 |
| OnAVOS [32] | ✓ | ✓ | | 86.1 | 84.9 | 85.5 | 0.08 |
| PReMVOS [18] | ✓ | ✓ | | 84.9 | 88.6 | 86.8 | 0.02 |
| STM-ResNet50 [20]$^{\dagger}$ | ✓ | | | 88.9 | 88.9 | 88.9 | 7.4 |
| STM-ResNet50-cycle (Ours)$^{\dagger}$ | ✓ | | ✓ | **89.2** | **90.4** | **89.8** | 2.0 |
| STM-ResNet18-cycle (Ours) | | | | 80.4 | 80.3 | 80.4 | **64.5** |
| STM-ResNet18-cycle (Ours) | | | ✓ | 81.3 | 81.1 | 81.2 | 17.7 |
| STM-ResNet50-cycle (Ours) | | | | 84.1 | 83.7 | 83.9 | 38.5 |
| STM-ResNet50-cycle (Ours) | | | ✓ | 84.1 | 83.8 | 84.0 | 11.5 |

Table 1: Comparison with state-of-the-art method on DAVIS16 and DAVIS17 set. "Extra data" indicates the method is pretrained with extra data with mask annotations. "-" indicates unavailable results. "OL" denotes online learning or update process. "GC" is short for gradient correction. "$\dagger$" denotes the pretraining is implemented from [1] with larger input size ($480 \times 854$) for inference.

with the configuration $\mathcal{X}_{t-1} = \{X_{t-1}\}, \mathcal{Y}_{t-1} = \{\widehat{Y}_{t-1}\}$, since this case is the most vulnerable to accumulative error propagation and hence, proper training with cyclic loss term can effectively relieve such a problem.

*4.3.3 Sensitivity analysis.*

Next, we evaluate how the hyperparameters in our algorithm affect the final results. In Figure 4, we show the performance-runtime trade-off w.r.t. the correction iteration time $N$. We find that the $\mathcal{J}\&\mathcal{F}$ score saturates when $N$ approaches 10; above which, the score improvement is somewhat marginal but at the expense of decreasing efficiency. Considering the trade-off between runtime and performance, we take $N = 10$ as the empirically optimal iteration number for gradient correction. Additionally, we also analyze the impact of correction rate $\alpha$ as shown in Figure 7. From this figure, we find the larger strength of correction usually results in better performance, but the overall performance vari-

| Method | Extra data | OL | GC | $\mathcal{J}_\mathcal{S}(\%)$ | $\mathcal{J}_\mathcal{U}(\%)$ | $\mathcal{F}_\mathcal{S}(\%)$ | $\mathcal{F}_\mathcal{U}(\%)$ | $\mathcal{G}(\%)$ | FPS |
|---|---|---|---|---|---|---|---|---|---|
| RVOS [30] | | | | 63.6 | 45.5 | 67.2 | 51.0 | 56.8 | 24 |
| S2S [35] | | | | 66.7 | 48.2 | 65.5 | 50.3 | 57.6 | 6 |
| FRTM [25] | | | | 68.6 | 58.4 | 71.3 | 64.5 | 65.7 | - |
| TVOS [38] | | | | 67.1 | 63.0 | 69.4 | 71.6 | 67.8 | - |
| RGMP [34] | ✓ | | | 59.5 | - | 45.2 | - | 53.8 | 7 |
| DMM-Net [37] | ✓ | | | 58.3 | 41.6 | 60.7 | 46.3 | 51.7 | 12 |
| AGSS-VOS [16] | ✓ | | | 71.3 | 65.5 | 75.2 | 73.1 | 71.3 | 12.5 |
| S2S [35] | | | ✓ | 71.0 | 55.5 | 70.0 | 61.2 | 64.4 | 0.06 |
| OSVOS [3] | ✓ | | ✓ | 59.8 | 54.2 | 60.5 | 60.7 | 58.8 | - |
| MaskTrack [22] | ✓ | | ✓ | 59.9 | 45.0 | 59.5 | 47.9 | 53.1 | 0.05 |
| OnAVOS [32] | ✓ | | ✓ | 60.1 | 46.6 | 62.7 | 51.4 | 55.2 | 0.05 |
| DMM-Net [37] | ✓ | | ✓ | 60.3 | 50.6 | 63.5 | 57.4 | 58.0 | - |
| STM-ResNet50 [20]$^\dagger$ | ✓ | | | 76.1 | 70.8 | 79.6 | 77.5 | 76.0 | 3.9 |
| STM-ResNet50-cycle (Ours)$^\dagger$ | ✓ | | ✓ | **77.8** | **73.3** | **81.5** | **80.1** | **78.2** | 0.9 |
| STM-ResNet18-cycle (Ours) | | | | 69.2 | 56.2 | 72.5 | 65.0 | 65.7 | **63** |
| STM-ResNet18-cycle (Ours) | | | ✓ | 70.4 | 58.2 | 73.9 | 67.2 | 67.5 | 15.2 |
| STM-ResNet50-cycle (Ours) | | | | 71.7 | 61.4 | 75.8 | 70.4 | 69.9 | 43 |
| STM-ResNet50-cycle (Ours) | | | ✓ | 72.6 | 63.0 | 76.7 | 72.3 | 71.2 | 13.8 |

Table 2: Comparison with state-of-the-art method on Youtube-VOS validation set. The subscript $\mathcal{S}$ and $\mathcal{U}$ denote the seen and unseen categories. $\mathcal{G}$ is the global mean. "-" indicates unavailable results. "OL" denotes online learning or update process. "GC" is short for gradient correction. "$\dagger$" denotes the pretraining is implemented from [1] with larger input size ($480 \times 854$) for inference.

| datasets | | DAVIS17 | | | DAVIS16 | | |
|---|---|---|---|---|---|---|---|
| $\mathcal{X}_{t-1}$ | $\mathcal{Y}_{t-1}$ | baseline | +cycle | $\Delta$ | baseline | +cycle | $\Delta$ |
| $\{X_1\}$ | $\{Y_1\}$ | 65.2 | 67.6 | +2.4 | 81.2 | 81.2 | +0.0 |
| $\{X_{t-1}\}$ | $\{\widehat{Y}_{t-1}\}$ | 56.8 | 61.2 | **+4.4** | 75.3 | 77.2 | **+1.9** |
| $\{X_1, X_{t-1}\}$ | $\{Y_1, \widehat{Y}_{t-1}\}$ | 67.3 | 69.2 | +1.9 | 82.3 | 83.8 | +1.5 |
| **MEM** | **MEM** | 69.7 | 71.7 | +2.0 | 82.3 | 83.9 | +1.6 |

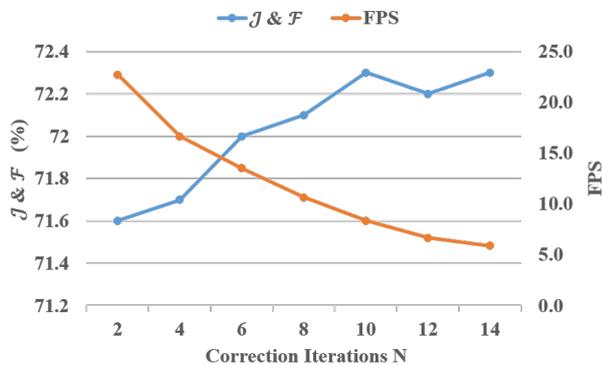Table 3: Experiments on improvement of $\mathcal{J}\&\mathcal{F}$ score with different reference set configuration.



Fig. 4: Performance-runtime trade-off with different iteration size $N$ on DAVIS17 validation.



Fig. 5: Results on DAVIS17 validation set with different amount of training data in Youtube-VOS for STM-cycle and baseline STM model.

ation is not sensitive to the change of correction rate $\alpha$, reflecting that our update scheme is robust and can accommodate variations to this parameter well, consequently, we set correction rate $\alpha = 180$ since the overall gain becomes saturated under this configuration.
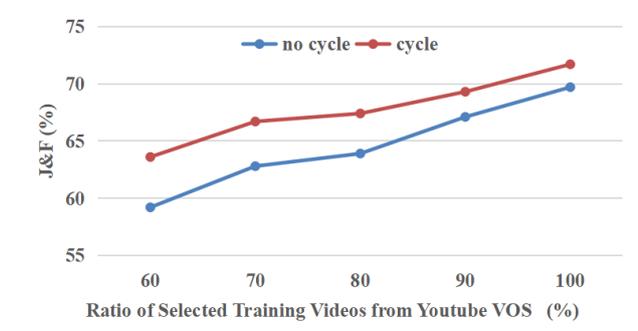
### 4.3.4 Effect of Anti-noise regularization

Together with the analysis of correction rate $\alpha$, we also investigate the effect of proposed anti-noise regularization term under different correction rate and loss weight $\lambda$. We find that the smooth term in Equation (7) can bring stable gains to overall segmentation quality except when the correction rate is small. When $\alpha$ is larger, the gain appears to be more obvious and finally con-

| | datasets | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&djk\mathcal{F}(\%)$ |
|---|---|---|---|---|
| baseline | | 67.6 | 71.7 | 69.7 |
| + cyclic | DAVIS17 | 68.7 | 74.7 | 71.7 |
| + GC | | 69.2 | 74.3 | 71.8 |
| + both | | **69.8** | **75.9** | **72.9** |
| baseline | | 82.8 | 81.7 | 82.3 |
| + cyclic | DAVIS16 | 84.1 | 83.7 | 83.9 |
| + GC | | 83.3 | 82.3 | 82.8 |
| + both | | **84.1** | **83.8** | **84.0** |

Table 4: Ablation study on the effectiveness of different component. "GC" is short for gradient correction.



Fig. 6: Qualitative comparison between predicted results from gradient correction module with and without anti-noise regularization on DAVIS17 test-dev. Left column: The original query frame. Medium column: Predicted mask from gradient correction module without anti-noise regularization. Right column: Predicted mask from gradient correction module with anti-noise regularization.
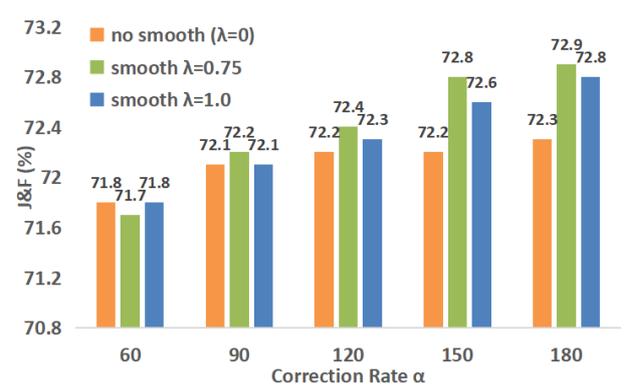


Fig. 7: Performance with different correction rate and smooth term for STM-cycle. The model is trained on the hybrid dataset of DAVIS2017 and Youtube-VOS. And the evaluation results are reported on DAVIS2017 validation set.

verges to a stable level of improvement. We think this is because larger correction rates produces more precise segmentation but also tends to result in more severe background artifacts, which can be alleviated by the smooth constraint. Finally, we set $\lambda = 0.75$ since we find the overall gain is not sensitive within proper interval of this loss weight.

In Figure 6, we further qualitatively analysis the impact of the regularization term during gradient correction. By comparison between the output masks, we can find we equipped with the smooth regularization, our gradient correction can suppress some subtle noise blocks in the background.

### 4.3.5 Results with Different Amount of Training Data

As discussed in Section 3.2, the cycle-consistency loss in our method implicitly applies data augmentation to train more generalized VOS model, thus should be more robust to scarcity of training data, especially for visual tasks that requires large amounts of synthesized

data [17] for pretraining like VOS. To demonstrate this discussion, we specifically design experiments to test the performance when the training data is gradually decreased. In detail, we start from our standard benchmark with hybrid training set of DAVIS and Youtube-VOS (%100 selected), then we gradually decrease the available training clips in Youtube-VOS (from %100 to %60) and evaluate the performance of trained model. The results are depicted in Figure 5. We can observe that when combined cyclic consistency loss, the performance of trained STM model degrades slower than the counterpart without cycle-consistency but with the same backbone, demonstrating that the cycle consistency can alleviate the problem of data scarcity to some extent.

### 4.3.6 Robustness to noise perturbations

In addition to accumulative predicted error, video object segmentation systems are also prone to noise perturbations on reference templates. In this section, we further investigate how the gradient correction process mitigates the effects from such noisy interference. To do this, we utilize STM network with the **MEM** strategy as [20] by dynamically appending a predicted mask and its frame into the reference set. However, in this case, the predicted masks to be appended are manually replaced by a noisy version. Technically, we try to perturb the inference process with four types of noises, of which two are regarded as natural noises and the other two are adversarial-type noises.

– **Low-quality**. We replace the predicted mask $\widehat{Y}_t$ from baseline model with lower segmentation quality on the same frame.

| datasets | | DAVIS17 | | | | DAVIS16 | | |
|---|---|---|---|---|---|---|---|---|
| Noise | +GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | +GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ |
| Low-quality | | 65.9 | 72.2 | 69.1 | | 65.9 | 72.2 | 69.1 |
| | ✓ | **66.9** | **73.3** | **70.1** | ✓ | **66.9** | **73.3** | **70.1** |
| Box-template | | 63.0 | 66.0 | 64.5 | | 63.0 | 66.0 | 64.5 |
| | ✓ | **68.7** | **74.9** | **71.8** | ✓ | **68.7** | **74.9** | **71.8** |

Table 5: Results of models affected by natural noisy masks on DAVIS17 and DAVIS16 validation set.

| setting | | White Box | | | | Black Box | | |
|---|---|---|---|---|---|---|---|---|
| Noise | +GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | +GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ |
| FGSM [7] | | 43.2 | 48.8 | 46.0 | | 39.0 | 46.4 | 42.7 |
| | ✓ | **52.7** | **59.4** | **56.0** | ✓ | **51.0** | **59.7** | **55.4** |
| MI-FGSM [5] | | 38.6 | 44.8 | 41.7 | | 30.9 | 36.8 | 33.9 |
| | ✓ | **50.6** | **58.0** | **54.3** | ✓ | **47.7** | **55.6** | **51.7** |

Table 6: Results of models affected by adversarial noise on DAVIS17 validation set.

| datasets | setting | | DAVIS17 | | | DAVIS16 | | |
|---|---|---|---|---|---|---|---|---|
| base model | +cycle | +GC | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ | $\mathcal{J}(\%)$ | $\mathcal{F}(\%)$ | $\mathcal{J}\&\mathcal{F}(\%)$ |
| STM [20] | | | 67.6 | 71.7 | 69.7 | 82.8 | 81.8 | 82.3 |
| | ✓ | | 68.7 | 74.7 | 71.7 | 84.1 | 83.5 | 83.8 |
| | ✓ | ✓ | **69.8** | **75.9** | **72.9** | **84.1** | **83.8** | **84.0** |
| KMN [27] | | | 67.5 | 72.3 | 69.9 | 81.5 | 80.3 | 80.9 |
| | ✓ | | 67.8 | 73.5 | 70.7 | 82.6 | 83.4 | 83.0 |
| | ✓ | ✓ | **69.1** | **75.3** | **72.2** | **82.8** | **83.7** | **83.3** |
| GCM [14] | | | 66.7 | 72.9 | 69.8 | 81.1 | 81.1 | 81.1 |
| | ✓ | | 67.5 | 73.3 | 70.4 | 83.3 | 83.2 | 83.3 |
| | ✓ | ✓ | **67.6** | **73.7** | **70.7** | **83.6** | **83.6** | **83.6** |
| AOT-T [36] | | | 76.5 | 81.9 | 79.2 | 86.5 | 88.4 | 87.5 |
| | ✓ | | 77.7 | 82.7 | 80.2 | 86.5 | 88.5 | 87.5 |
| | ✓ | ✓ | **77.9** | **83.6** | **80.8** | **86.6** | **88.5** | **87.6** |

Table 7: Results of cyclic mechanism with different baseline models on DAVIS17 and DAVIS16 validation sets. The baseline results of STM, KMN and GCM are from our own implementation. The results of AOT-T are obtained from the official public code.

- **Box-template**. We replace the predicted mask $\widehat{Y}_t$ with a coarse level groundtruth mask where all pixels in the bounding box of objects are set to be 1.
- **FGSM** [7]. We take the classical adversarial attack method to generate interference noise on intermediate reference masks, where the loss term in Equation (4) on all frames in a video is maximized under a given pixel-wise changing constraint $\epsilon \leq 20$.
- **MI-FGSM** [5]. We further added a harsher adversarial attack method to test the robustness of our model. The MI-FSGM method generates noise towards the same objective and constraint as FSGM but updates the reference mask according to an iterative manner with momentum.

For the noise types generated by adversarial attack – FGSM [7] and MI-FGSM [5], we follow the common protocol in adversarial attack and defence by testing the results under both black box and white box settings. In white box setting, we take the trained STM-ResNet50-cycle model to generate noise and the attack itself, while in black box setting, we leverage the trained

STM-ResNet18-cycle model to generate noise and attack the one with the vanilla ResNet50 backbone. For each scheme, we conduct another experiment with gradient correction on the noisy masks before appending to the memory as the control group. The results are reported in Table 5 and Table 6. From Table 5, we see the gradient correction is helpful for both low-quality and box-template reference conditions. The improvement is much more obvious for the case of box-template, which indicates that the impact of gradient correction is greater when the intermediate reference mask is coarser but properly covers the object area. Meanwhile, from Table 6, we could see that although both attack methods degrade the segmentation performance to a large extent, but by properly utilizing the gradient correction, we can alleviate the effect from adversarial attacks.

### 4.3.7 Extension to other memory-based methods.

Our cyclic mechanism is easy to implement and can be naturally extend to other segmentation pipelines. To
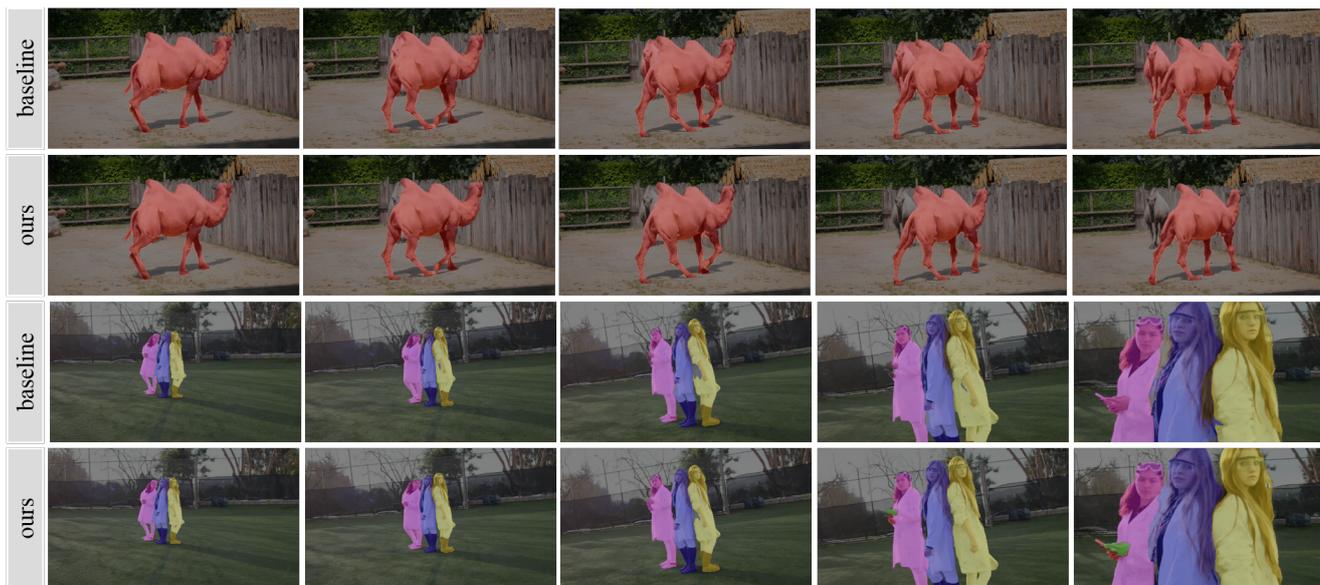
Fig. 8: Qualitative results shows the improvement of cyclic training over the baseline in DAVIS17 validation.
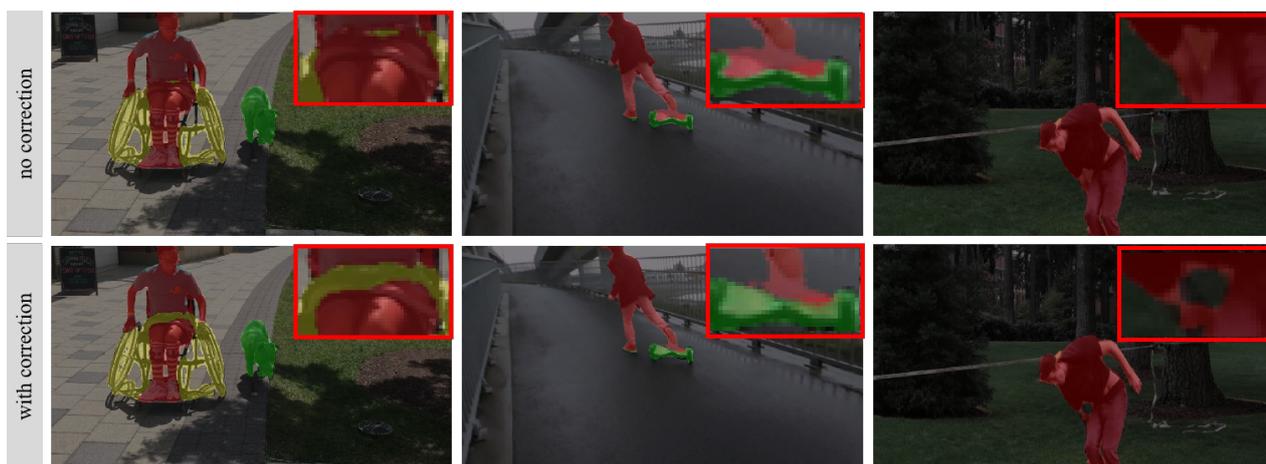


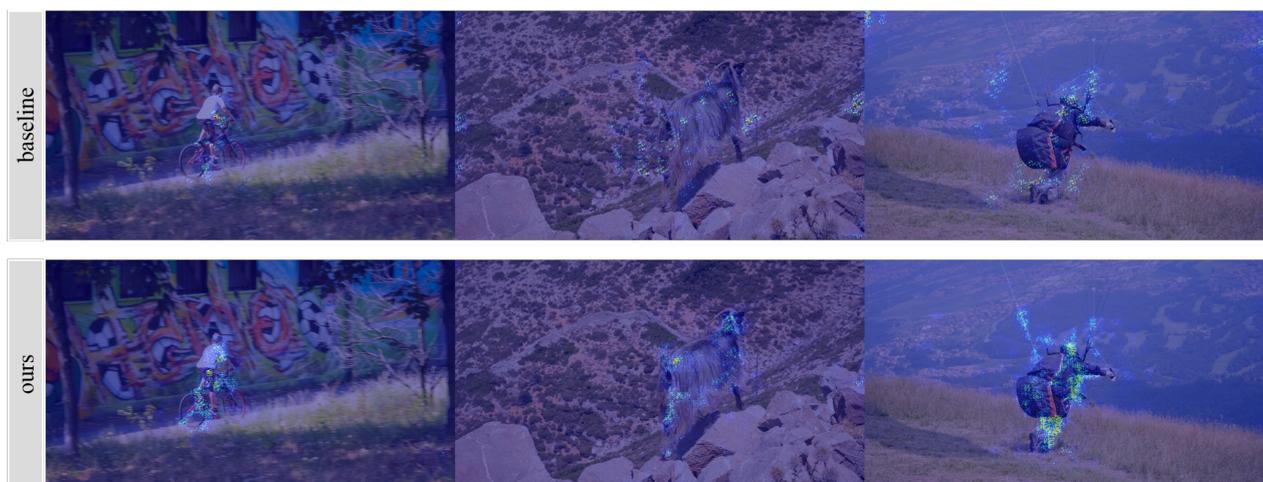Fig. 9: The visual effect of gradient correction module on DAVIS17 test-dev set



Fig. 10: Cycle-ERF of frames w.r.t. the initial reference object masks in DAVIS17 validation.

Fig. 11: Cycle-ERF comparison on STM, GCM and KMN on DAVIS16 validation set.

demonstrate the generalization of our design, we further extend another two famous baseline methods, KMN [27] GCM [14] and AOT [36] to a cyclic version and measure the performance gain. Corresponding results are reported in Table 7, we can observe that combining with cyclic mechanism can boost the segmentation quality for all the baseline models on both datasets. In detail, we find the improvement of gradient correction on GCM is the least obvious, this can be due to the global property of GCM, where the memory squeezes the spatial domain, and loses the fine-grain information of original objects, which prevents the gradient from propagating to detailed spatial location. In contrast, the gain of gradient correction on KMN is most salient, since KMN itself will highlight the focusing area in the query frame, therefore the backward gradient related to the objects will be preserved. This is also reflected in the cycle-ERF analysis in section 4.4.

## 4.4 Qualitative analysis

**Segmentation results.** In Figure 8, we show some segmentation results using the STM model trained with and without our cycle scheme. From comparison on the first sequences, we observe that the cyclic mechanism suppresses the accumulative error from problematic reference masks. From the second video, we see the cyclic model can depicts the boundary between foreground

objects more precisely, which is consistent with the quantitative results. Further, our method can successfully segment some challenging small objects (caught by the left woman's hand). In Figure 9, we show the visual effect from the gradient correction module, it is clear to see that gradient correction can help update and reconstruct the losses on some detail.

**Cycle-ERF analysis.** We further analyze the cycle-ERF defined as Equation (12) on different approaches. We take the initial mask as the objects to be predicted and take a random intermediate frame and an empty mask as reference. Figure 10 visualizes the cycle-ERFs of some samples output from baseline STM and STM-cycle model. Compared with baseline, our cyclic training scheme helps the network concentrate more on the foreground objects with stronger responses. This indicates that our model learns more robust object-specific correspondence. It is also interesting to see that for STM network, only a small part of the objects is crucial for reconstructing the same objects that were in the initial frames as the receptive field focuses on the outline or skeleton of the objects. This can be used to explain the greater improvement of contour accuracy using our method, and also provide cues on the extraction from reference masks.

In Figure 11, we take the cycle-ERF as a tool to analysis the focusing area of different baseline models, STM [20], GCM [14] and KMN [27] on DAVIS16. By comparison, we find STM shows overall stronger re-

Fig. 12: Failure cases of STM-cycle on DAVIS17 test-dev.

sponse on focusing area around the foreground object. In contrast, GCM highlight more background context around the object, since the memorial mechanism in this method always squeeze the context into the cache and focus more on global relationship. On the other hand, cycle-ERF from KMN yields response focusing on more specific and small part of objects and less intensity in background or context, this is in consistency with the design insight behind this method, where a dynamic gaussian kernel is applied to suppress the interaction between less related areas. These comparisons with Cycle-ERF manifest that such visualization methods can be a helpful tool to provide interpretability of existing models.

**Investigation of Failure Cases.** Figure 12 shows some failure cases of our method, although combined with cyclic loss and gradient correction, the network can not handle extremely narrow and small objects (e.g. the brassie in the man's hand in the second row), meanwhile, as shown in the first row, our method can suffer from cases where the specified objects are severely occluded by obstacles in the foreground.

## 5 Conclusion

This paper incorporates the cycle mechanism with semi-supervised video segmentation network to mitigate the error propagation problem in current approaches. When combined with an efficient and flexible baseline, the proposed cyclic loss and gradient correction module achieve competitive performance-runtime trade-off on three challenging benchmarks. Detailed analysis are further conducted to demonstrate the generality and robustness of such cyclic design. Further explanations can be drawn from a new perspective of cycle-ERF.

## References

1. Pretraining code of space-time-memory network on coco for video object segmentation. https://github.com/haochenheheda/Training-Code-of-STM
2. Bansal, A., Ma, S., Ramanan, D., Sheikh, Y.: Recycle-gan: Unsupervised video retargeting. In: European Conference on Computer Vision (ECCV) (2018)
3. Caelles, S., Maninis, K.K., Pont-Tuset, J., Leal-Taixé, L., Cremers, D., Van Gool, L.: One-shot video object segmentation. In: Computer Vision and Pattern Recognition (CVPR) (2017)
4. Carl, V., Abhinav, S., Alireza, F., Sergio, G., Kevin, M.: Tracking emerges by colorizing videos. European Conference on Computer Vision (2018)
5. Dong, Y., Liao, F., Pang, T., Su, H., Zhu, J., Hu, X., Li, J.: Boosting adversarial attacks with momentum. In: Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR), pp. 9185–9193 (2018)
6. Everingham, M., Eslami, S.M.A., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes challenge: A retrospective. International Journal of Computer Vision (IJCV) **111**(1), 98–136 (2015)
7. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
8. He, K., Georgia, G., Piotr, D., Ross, G.: Mask r-cnn. In: International Conference on Computer Vision (ICCV) (2018)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
10. Jabri, A., Owens, A., Efros, A.A.: Space-time correspondence as a contrastive random walk. Advances in Neural Information Processing Systems (2020)
11. Johnander, J., Danelljan, M., Brissman, E., Khan, F.S., Felsberg, M.: A generative appearance model for end-to-end video object segmentation. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 8953–8962 (2019)
12. Khoreva, A., Benenson, R., Ilg, E., Brox, T., Schiele, B.: Lucid data dreaming for video object segmentation. International Journal of Computer Vision (IJCV) **127**, 1175–1197 (2019)
13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representation (ICLR) (2014)
14. Li, Y., Shen, Z., Shan, Y.: Fast video object segmentation using the global context module. In: European Conference on Computer Vision (ECCV), pp. 735–750 (2020)

15. Li, Y., Xu, N., Jinlong, P., See, J., Weiyao, L.: Delving into the cyclic mechanism in semi-supervised video object segmentation. In: Neural Information Processing System (NeurIPS) (2020)

16. Lin, H., Qi, X., Jia, J.: Agss-vos: Attention guided single-shot video object segmentation. In: The IEEE International Conference on Computer Vision (ICCV) (2019)

17. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European Conference on Computer Vision (ECCV). Zürich (2014)

18. Luiten, J., Voigtlaender, P., Leibe, B.: Premvos: Proposal-generation, refinement and merging for video object segmentation. In: Asian Conference on Computer Vision (ACCV) (2018)

19. Meister, S., Hur, J., Roth, S.: UnFlow: Unsupervised learning of optical flow with a bidirectional census loss. In: AAAI. New Orleans, Louisiana (2018)

20. Oh, S.W., Lee, J.Y., Xu, N., Kim, S.J.: Video object segmentation using space-time memory networks. In: The IEEE International Conference on Computer Vision (ICCV) (2019)

21. Peng, J., Wang, C., Wan, F., Wu, Y., Wang, Y., Tai, Y., Wang, C., Li, J., Huang, F., Fu, Y.: Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking. In: The European Conference on Computer Vision (ECCV) (2020)

22. Perazzi, F., Khoreva, A., Benenson, R., Schiele, B., Sorkine-Hornung, A.: Learning video object segmentation from static images. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

23. Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: Computer Vision and Pattern Recognition (CVPR) (2016)

24. Pont-Tuset, J., Perazzi, F., Caelles, S., Arbeláez, P., Sorkine-Hornung, A., Van Gool, L.: The 2017 davis challenge on video object segmentation. arXiv:1704.00675 (2017)

25. Robinson, A., Lawin, F.J., Danelljan, M., Khan, F.S., Felsberg, M.: Learning fast and robust target models for video object segmentation. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)

26. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision (IJCV) **115**(3), 211–252 (2015)

27. Seong, H., Hyun, J., Kim, E.: Kernelized memory network for video object segmentation. In: European Conference on Computer Vision (ECCV), pp. 629–645 (2020)

28. Shi, J., Yan, Q., Xu, L., Jia, J.: Hierarchical image saliency detection on extended cssd. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **38**(4), 717–729 (2016)

29. Tinghui, Z., Philipp, K., Mathieu, A., Qixing, H., Alexei, A.E.: Learning dense correspondence via 3d-guided cycle consistency. The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

30. Ventura, C., Bellver, M., Girbau, A., Salvador, A., Marques, F., Giro-i Nieto, X.: Rvos: End-to-end recurrent network for video object segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

31. Voigtlaender, P., Chai, Y., Schroff, F., Adam, H., Leibe, B., Chen, L.C.: Feelvos: Fast end-to-end embedding learning for video object segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

32. Voigtlaender, P., Leibe, B.: Online adaptation of convolutional neural networks for video object segmentation. In: British Machine Vision Conference (BMVC) (2017)

33. Wang, X., Jabri, A., Efros, A.A.: Learning correspondence from the cycle-consistency of time. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

34. Wug Oh, S., Lee, J.Y., Sunkavalli, K., Joo Kim, S.: Fast video object segmentation by reference-guided mask propagation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

35. Xu, N., Yang, L., Fan, Y., Yang, J., Yue, D., Liang, Y., Price, B., Cohen, S., Huang, T.: Youtube-vos: Sequence-to-sequence video object segmentation. In: The European Conference on Computer Vision (ECCV) (2018)

36. Yang, Z., Wei, Y., Yang, Y.: Associating objects with transformers for video object segmentation. In: Advances in Neural Information Processing Systems, vol. 34, pp. 2491–2502 (2021)

37. Zeng, X., Liao, R., Gu, L., Xiong, Y., Fidler, S., Urtasun, R.: Dmm-net: Differentiable mask-matching network for video object segmentation. In: The IEEE International Conference on Computer Vision (ICCV) (2019)

38. Zhang, Y., Wu, Z., Peng, H., Lin, S.: A transductive approach for video object segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2020)

39. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: IEEE International Conference on Computer Vision (ICCV) (2017)