# Efficient Long-Short Temporal Attention Network for Unsupervised Video Object Segmentation

Ping Li[a,b], Yu Zhang[a], Li Yuan[c], Huaxin Xiao[d], Binbin Lin[e,*], Xianghua Xu[a]

[a]*School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou, China*
[b]*Guangdong Laboratory of Artificial Intelligence and Digital Economy (SZ), Shenzhen, China*
[c]*School of Electrical and Computer Engineering, Peking University, Beijing, China*
[d]*College of Systems Engineering, National University of Defense Technology, Changsha, China*
[e]*College of Computer Science, Zhejiang University, Hangzhou, China*

## Abstract

Unsupervised Video Object Segmentation (VOS) aims at identifying the contours of primary foreground objects in videos without any prior knowledge. However, previous methods do not fully use spatial-temporal context and fail to tackle this challenging task in real-time. This motivates us to develop an efficient *Long-Short Temporal Attention* network (termed **LSTA**) for unsupervised VOS task from a holistic view. Specifically, LSTA consists of two dominant modules, i.e., Long Temporal Memory and Short Temporal Attention. The former captures the long-term global pixel relations of the past frames and the current frame, which models constantly present objects by encoding appearance pattern. Meanwhile, the latter reveals the short-term local pixel relations of one nearby frame and the current frame, which models moving objects by encoding motion pattern. To speedup the inference, the efficient projection and the locality-based sliding window are adopted to achieve nearly linear time complexity for the two light modules, respectively. Extensive empirical studies on several benchmarks have demonstrated promising performances of the proposed method with high efficiency.

*Keywords:* Unsupervised video object segmentation, long temporal memory, short temporal attention, efficient projection

## 1. Introduction

Video Object Segmentation (VOS) task is to localize and segment primary objects in videos, i.e., yielding accurate contours of objects. As a fundamental video processing technique, VOS has found widespread applications, e.g., video editing [35], autonomous driving, and surveillance environment [53], which are highly demanding in real-time processing. Generally, VOS methods are divided into two categories, i.e., *semi-supervised* VOS (*a.k.a.*, one-shot VOS) [12] which utilizes given object mask of the first frame, and *unsupervised* VOS (*a.k.a.*, zero-shot VOS) [20] for which arbitrary prior knowledge is unavailable during inference. This work concentrates on the more challenging unsupervised VOS, which faces **two** considerably critical problems: 1) how to find primary objects in video frames; 2) how to speedup object segmentation inference.

For the first problem, the common insight is to consider salient objects, moving objects, and constantly present objects across video frames. While salient objects attract the visual attention from human eyes, fast moving and drastic deformations may yield objects with small appearance size, leading to less saliency. To model moving objects, someone [9] adopt optical flow technique to capture motion cues, but it is still difficult to discriminate moving objects from dynamic background and usually fails to identify objects in static scenes. From a holistic view, a natural idea is to observe whether there exist constantly present objects in the past frames, and then search objects with similar appearance in the current frame. This idea has been proved effective [20, 40] by using dot product attention to encode pixel-wise dense correlations of past

---

*Corresponding author
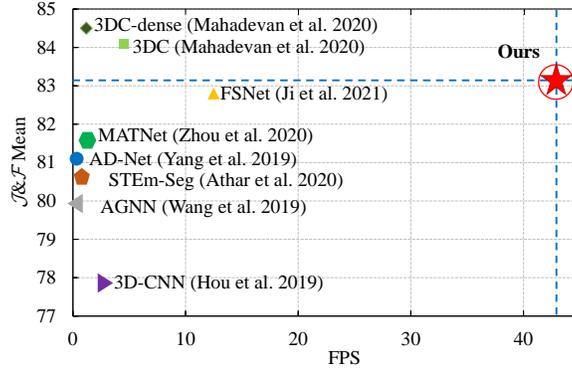*Email address:* binbinlin@zju.edu.cn (Binbin Lin)

Figure 1: Overall efficiency comparison of several SOTA unsupervised VOS methods without object prior on DAVIS2016 validation set.

frames. However, when partial area of objects are occluded, it adds much difficulty in identifying similar objects due to its strong reliance on appearance. To address these limitations, we model *constantly present objects* and *moving objects* at the same time, by utilizing motion cues and temporal consistency of objects in past frames from both *full-frame* (all pixels in a frame) and *partial-frame* (one frame is separated into many small patches) perspectives. Encoding full-frame pixel correlations facilitates tackling object deformation by modeling appearance pattern, while encoding partial-frame pixel correlations benefits handling object occlusion by modeling pixel movements in the local region of frame.

For the second problem, it still remains an open issue to be explored in unsupervised VOS without any object prior knowledge. Existing models cannot be deployed in real-time applications due to their low inference speed caused by using optical flow [55, 9] or 3D Convolutional Neural Networks (CNN) [22, 1], as illustrated in Fig. 1. Accordingly, we explore the way of accelerating inference from both full-frame and partial-frame perspectives to identify objects efficiently. As is well known, the time cost of directly encoding full-frame pixel correlation increases squarely with the number of pixels, which limits its applicability. Someone proposed channel-wise attention [17] to capture the global context of past frames for fast semi-supervised VOS, but it is unable to preserve per-pixel correlations, thus deteriorating performance. Therefore, inspired by the random projection on feature map for computing efficient attention relation with nearly linear complexity [11], we propose to adopt an efficient projection skill to reveal channel-wise correlation for unsupervised VOS by doing random projection on feature maps derived from CNNs. This projection can achieve the similarity distribution approximation of frames, such that the pixel-wise similarity between the past frames and the current frame can be well preserved in the embedding space. So it is considerably beneficial for discriminating constantly present objects. Meanwhile, since the number of channel $c$ is far less than that of pixel $n$ in a feature map, i.e., $c \ll n$, the time cost of encoding inherent relations among past frames is reduced from square complexity to linear level, e.g., $\mathcal{O}(n^2 c) \rightarrow \mathcal{O}(nc^2)$. On the other hand, the locality-based sliding window strategy is employed to partition one full frame to many overlapped patches with size of $k \times k$ ($k \ll n$), i.e., partial frames. This helps to model the local patterns of objects, such as edges, lines, and textures. By this means, encoding partial-frame pixel correlation requires linear time complexity, i.e., $\mathcal{O}(nck^2) = \mathcal{O}(nc) = \mathcal{O}(n)$, much less than directly encoding full-frame correlation, i.e., $\mathcal{O}(n^2 c)$.

Therefore, we propose an end-to-end real-time unsupervised VOS framework, named **L**ong-**S**hort **T**emporal **A**ttention network (**LSTA**), to strike a good balance between performance and speed. This framework mainly includes two fast modules, i.e., *Long Temporal Memory* (LTM) and *Short Temporal Attention* (STA). LTM enables encoding long-term full-frame pixel spatiotemporal dependency between the past frames and the current frame, which facilitates identifying constantly present objects. Simultaneously, STA enables capturing short-term partial-frame pixel spatiotemporal relations between one nearby frame and the current frame, which benefits finding moving objects. As a matter of fact, the two modules cooperatively work together to find primary objects by modeling both long-range and short-range spatiotemporal coherence of frames. Meanwhile, it paves the way for discriminating objects from complex background and thus alleviates the object deformation or occlusion problem. More importantly, we apply our proposed efficient projection to LTM and the locality-based sliding window to STA, respectively, for greatly reducing the time complexity. Thus, both LTM and STA can be implemented at linear time complexity, making our LSTA framework very

2

efficient. To examine its performance, we have conducted comprehensive experiments on several benchmark databases, i.e., DAVIS2016[26], DAVIS2017[27], YouTube-Objects[28], and FBMS[24]. Empirical studies demonstrate that our method exhibits promising segmentation performances at a fast speed, e.g., 42.8 fps on 480p resolution videos from DAVIS2016.

Our main contributions are highlighted in the following:

- We propose an end-to-end real-time unsupervised VOS framework, called Long-Short Temporal Attention network (LSTA), which enjoys satisfying segmentation accuracy with high inference efficiency.

- The Long Temporal Memory (LTM) module and the Short Temporal Attention (STA) module are developed to encode both global and local spatiotemporal pixel-wise relations among frames. Hence, constantly present and moving objects can be readily found, and the object deformation or occlusion problem can be alleviated.

- LTM module and STA module can both achieve the nearly linear time complexity, by respectively adopting the efficient projection and the locality-based sliding window strategy on feature maps.

- Performance comparisons and extensive ablation studies have justified the realtime segmentation ability with high precision by our method on several benchmarks.

The rest of this paper is organized as follows. Section 2 reviews closely related works and Section 3 introduces the newly developed LSTA framework. After that, we report both quantitative and qualitative experimental results to verify the efficacy of the proposed method in Section 4. Finally, we conclude this work in Section 5.

## 2. Related Work

This section makes a brief summary of closely related VOS methods, including unsupervised, semi-supervised, and fast scenarios. Note that *unsupervised* and *semi-supervised* terms are indicated by whether using the first frame mask during inference. This is a bit different from the traditional machine learning paradigm. For a thorough survey on video segmentation using deep learning techniques, please refer to [56].

### 2.1. Unsupervised VOS

Unsupervised VOS methods have no prior on the first frame mask for inference, making it fairly challenging. Usually, existing methods attempt to find primary objects by considering temporal motion [55] or object saliency [33]. Here, we review two primary kinds of unsupervised methods, i.e., attention-based, and optical flow-based, which adopt 2D convolutions.

**Attention-based** methods [20][40][48] find objects using appearance feature derived from video frames. For example, COSNet (Co-Attention Siamese Networks) [20] captures global per-pixel correlation and scene context by using co-attention mechanism on visual features of different frames, which helps to find constantly present objects. But COSNet models spatiotemporal relation between only two nearby frames during inference, which easily causes error accumulation by iterative updates and fails to well capture long-range context of frames. To overcome this drawback, AGNN (Attentive Graph Neural Networks) [40] builds fully-connected graph, where a node is the frame feature and an edge stands for the relation of pair-wise features. However, AGNN largely relies on object appearance similarity, and performs poorly when partial objects are occluded. Both COSNet and AGNN utilize dot product attention that requires intensive computations, consequently preventing them from being widely deployed. While attention-based methods focus more on object appearance, partial background areas sharing similar appearance with primary objects will be treated as objects by mistake. To handle this shortcoming, AD-Net (Anchor Diffusion Network) [48] adopts instance pruning as a postprocessing to filter out some noisy objects via object detection. In addition, AGS (Attention-Guided object Segmentation) [41] computes visual attention using eye tracking data, and obtains the coarse object location through dynamic visual attention prediction. Nevertheless, AGS employs ConvLSTM (Convolutional Long Short-Term Memory) to model temporal relations, which fails to fully model long-range spatiotemporal context of frames. And a variant of ConvLSTM named RNN-Conv [53] aggregates the temporal and the spatial information, such that the model can discover important objects in video.

3

**Optical flow-based** methods [57] capture motion cues from optical flow feature as the compensation of appearance feature. The early work Segflow [5] unifies CNN and optical flow prediction network to predict object mask and optical flow simultaneously, which obtains motion cues in an end-to-end manner; another work [42] produces a spatiotemporal edge map by combining static edge probability and optical flow gradient magnitude. But they fail to fully use object appearance features, leading to inferior performance. Thereafter, some works [55][9] concentrate on how to derive and then fuse both the motion feature and the appearance feature. For instance, MATNet (Motion-Attentive Transition Network) [55] employs dot product attention to fuse motion and appearance features; RTNet (Reciprocal Transformation Network) [29] mutually evolves the two modalities such that the intra-frame contrast, the motion cues, and temporal coherence of recurring objects are holistically considered; TransportNet [51] employs the Wasserstein distance compute the global optimal flows to transport the features in one modality to the other, and formulates the motion-appearance alignment as an instance of optimal structure matching. But they require heavy computations, and to reduce time cost, FSNet (Full-duplex Strategy Network) [9] makes feature fusion by channel-wise attention, but distills out effective squeezed cues from feature, readily overlooking appearance details. Except for Segflow, the other methods require additional optical flow features, which are not end-to-end and also time-consuming. Besides, optical flow feature mainly encodes the temporal relation between only two nearby frames, failing to model the long-range relation. This usually makes the model perform not well, when drastic changes happen to objects in a long video.

### 2.2. Semi-supervised VOS

Semi-supervised VOS methods [7, 49] aim to capture objects in video given the first frame mask, which is class-agnostic. Previous methods can be roughly separated into three groups, i.e., *online learning-based*, *attention-based*, *detection-based*.

**Online learning-based** methods [31] employ the first frame and its mask to update model parameters during inference, which can adapt to the videos containing various objects in different categories. For example, Sun *et al.* [35] utilize reinforcement learning to select optimal adaptation areas for each frame, and make the model take optimal actions to adjust the region of interest inferred from the previous frame for online model updating; Lu *et al.* [19] perform the memory updating by storing and recalling target information from the external memory. However, the model updating is very slow, resulting in low segmentation speed.

**Attention-based** methods [25] treat the first frame mask as the model input to provide object prior during inference. They encode pairwise pixel dependency among video frames by attention mechanism, which helps to capture objects existing for a long time in the past frames. For example, MUNet (Motion Uncertainty-aware Net) [34] designs a motion-aware spatial attention module to fuse the appearance features and the motion uncertainty-aware feature. The drawback is the high computational cost of computing attention scores, i.e., pairwise pixel similarity.

**Detection-based** methods [46] always use object detection model to obtain object proposals in each frame, whose feature representations are propagated along the temporal dimension, and generate object masks in line with appearance similarity. The shortcoming is the quality of object proposals will heavily affect the mask quality, and they cannot be trained in an end-to-end way, leading to sub-optimal results.

### 2.3. Fast VOS

Most existing fast VOS methods [31][44][39][50] belong to semi-supervised paradigm, and they strive to efficiently extract discriminant features from video frames. For example, Robinson *et al.*[31] propose FRTM (Fast and Robust Target Models) that only updates partial model parameters to speedup inference for online learning-based methods. For attention-based methods, Li *et al.*[17] use channel-wise similarity to substitute pixel-wise similarity for capturing the global context among frames. Although this substitution can reduce the time complexity, the performance is still far from that of pixel-wise methods such as STM (Space-Time Memory Networks) [25]. Moreover, Swiftnet [39] uses sparse features, i.e., only computing the similarity of those more informative pixels, to reduce dense pixel computations in attention-based methods.

Real-time unsupervised VOS task still remains less explored in existing works. The one relevant work is WCS-Net (Weighted Correlation Siamese Network) [52] that borrows eye gaze estimation model to provide coarse object location, which is fed into a light segmentation model to obtain object masks of frames. While this approach achieves relatively high inference speed, it does require additional model to yield object prior information as pre-processing step. The other one is Dynamic Atrous Spatial Pyramid Pooling (ASPP) [53],
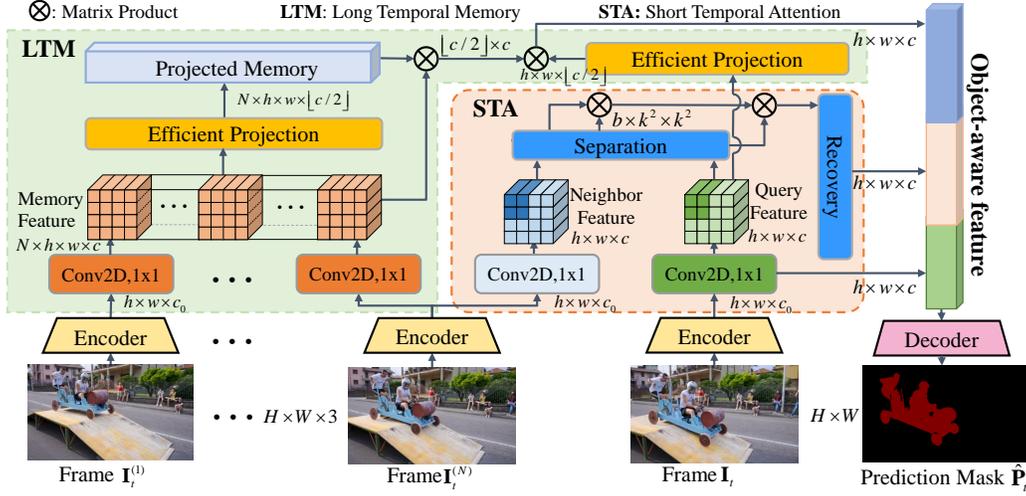
Figure 2: The framework of LSTA model for unsupervised VOS. It is composed of Encoder, LTM, STA, and Decoder. Note that, for the feature map with $c$ channels in LTM and STA, $h$ and $w$ denote height and width, respectively; for STA, $b$ is the number of local patches with size $k$ in one frame after passing the separation layer.

which adopts a dynamic selection mechanism in ASPP, and the dilated convolutional kernels adaptively select appropriate features by the channel attention mechanism. This still requires large computations with an additional RNN-Conv module. Luckily, our LSTA approach can achieve a good balance between segmentation accuracy and inference speed without any object prior, and can be trained in an end-to-end manner.

In addition, referring VOS has recently received more attention from the research field, such as Liang *et al.* [18] explore both local and global temporal context by an improved Transformer model to query the video with the language expression in an efficient manner. Also, the spatio-temporal context is important for weakly-supervised video grounding [15], which localizes the aligned visual tube corresponding to a language query. Very recently, Ji *et al.* [10] explored the Segment Anything Model (SAM) model on a variety of image segmentation tasks, such as agriculture, remote sensing, and healthcare, which may shed some light on the future research of unsupervised VOS task.

## 3. Our LSTA Method

To efficiently identify primary objects, we develop a real-time end-to-end unsupervised VOS approach, i.e., LSTA, by respecting the spatiotemporal coherence structure in frame data space. First of all, we briefly describe the problem formulation. Then, the main components including Encoder, Long Temporal Memory (LTM) block, Short Temporal Attention (STA) block, and Decoder, in the LSTA framework as illustrated by Fig. 2, will be elaborated.

### 3.1. Problem Formulation

Given a video sequence with $T$ frames, i.e., $\mathcal{V} = \{\mathbf{I}_t \in \mathbb{R}^{H \times W \times 3} | t = 1, 2, \ldots, T\}$, where $\mathbf{I}_t$ denotes the $t$-th RGB frame with width $W$, height $H$, and three channels. There may exist one or more than one objects in each frame, but no prior knowledge about objects are available. Unsupervised VOS aims to predict the pixel-wise object mask without specifying object in the first frame. For a video sequence $\mathcal{V}$ with one object, the ground-truth mask sequence is $\mathcal{P} = \{\mathbf{P}_t \in \{0, 1\}^{H \times W} | t = 1, 2, ..., T\}$ and the predicted mask sequence is $\hat{\mathcal{P}} = \{\hat{\mathbf{P}}_t \in \{0, 1\}^{H \times W} | t = 1, 2, ..., T\}$. The frame mask is a matrix with binary entries, where '0' means background pixel and '1' means primary object pixel. During training, the model uses RGB frames and their ground-truth masks as input, while only RGB frames are available during inference. Note that LSTA does not use all past frames but averagely divides them into $N$ bins, from each of which one frame is randomly selected. For the current frame $\mathbf{I}_t$, i.e., query frame, its past frame set is denoted as $\mathcal{I}_t = \{\mathbf{I}_t^{(1)}, \mathbf{I}_t^{(2)}, \ldots, \mathbf{I}_t^{(N)}\}$.

$\mathbf{F}_t^{(1)} \downarrow h \times w \times c_0$     $\mathbf{F}_t^{(N)} \downarrow h \times w \times c_0$     $\mathbf{F}_t \downarrow h \times w \times c_0$

$\phi$     ···     $\phi$     $\psi$

Reshape     Reshape     Reshape

$\tilde{\mathbf{F}}_t^{(1)}$     $\tilde{\mathbf{F}}_t^{(N)} \, hw \times c$     $\mathbf{Q}_t \, hw \times c$

$hw \times c$   Concatenate

$\mathbf{M}_t \, Nhw \times c$

Efficient Projection

$\mathbf{M}_t' \, Nhw \times \lfloor c/2 \rfloor$     $\mathbf{Q}_t' \, hw \times \lfloor c/2 \rfloor$

Norm

$\mathbf{D}_t \, hw \times c$

$\mathbf{M}_t$     $\hat{\mathbf{A}}_t$

$\lfloor c/2 \rfloor \times c$     $\mathbf{G}_t \, hw \times c$

$\mathbf{G}_t'$

Reshape   $hw \times c$

$\hat{\mathbf{G}}_t \downarrow h \times w \times c$     ⊘ : Element-wise Division
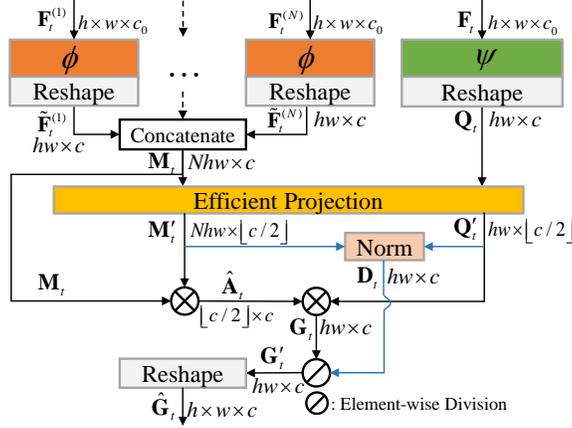
Figure 3: The Long Temporal Memory (LTM) block in LSTA. This block mainly consists of the convolution layer and orthogonal random projection, resulting in projected feature maps. Note that $\phi(\cdot)$ and $\psi(\cdot)$ are convolution layers acting as linear projection.

As illustrated in Fig. 2, Encoder adopts DeepLab v3+[3] without the last convolution layer, pre-trained on MS COCO database, and it is used to derive object-aware appearance feature from RGB frame $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$, where $H$ denotes height and $W$ denotes width. LTM models long-term full-frame pixel spatiotemporal relation between the past $N$ frames (memory) and the current frame (query) at time step $t$ using channel-wise attention with the efficient projection, which facilitates capturing constantly present objects. STA adopts the locality-based sliding window strategy in the separation layer and attention mechanism on the appearance features of the nearby frame $\mathbf{I}_t^{(N)}$ and the current frame $\mathbf{I}_t$, which helps to model the pattern of those moving objects. Decoder that consists of convolution layers, anisotropic convolution block [14], and bilinear up-sampling, is used for aggregating features derived from Encoder, LTM, and STA, resulting in object-aware feature representation for computing prediction mask $\hat{\mathbf{P}}_t \in \mathbb{R}^{H \times W}$ of the $t$-th frame. The recovery layer is used to reshape the feature map with the size of $b \times k^2 \times c$ to $h \times w \times c$.

### 3.2. Encoder

To encode the appearance property of video frames, we use the DeepLab v3+ [3] model pre-trained on MS COCO database as Encoder, and the last convolution layer is abandoned. As well known, DeepLab v3+ is a typical semantic segmentation model with ResNet101 as its backbone, and the pre-trained model can discriminate a large number of semantic classes of objects.

All frames in video sequence $\mathcal{V}$ are fed into Encoder to derive its appearance feature map. For the $t$-th frame and its past $N$ frames, we have

$$\{\mathbf{F}_t^{(1)}, \cdots, \mathbf{F}_t^{(N)}, \mathbf{F}_t\} = \{\Phi(\mathbf{I}_t^{(1)}), \cdots, \Phi(\mathbf{I}_t^{(N)}), \Phi(\mathbf{I}_t)\} \in \mathbb{R}^{h \times w \times c_0}, \tag{1}$$

where the function $\Phi(\cdot)$ denotes Encoder which projects RGB frame into feature map $\mathbf{F}_t$ with $c_0$ channels ($c_0$ is 256); the height is $h = \frac{H}{4}$ and the width is $w = \frac{W}{4}$. Since Encoder adopts pre-trained semantic segmentation model, it is able to capture the intrinsic appearance structure of common foreground objects in video frames. This is beneficial for finding those primary objects for segmentation.

### 3.3. Long Temporal Memory (LTM)

To identify those constantly present objects in video, LTM block, as illustrated in Fig. 3, employs appearance features of the past frames and the current frame to encode the full-frame pixel spatiotemporal dependency in terms of appearance similarity. This not only helps the model to readily find out those objects with similar appearance in the long-range frame context, such that constantly present objects in the current frame receive more attention, but also makes the model robust to object deformation.

To encode the spatiotemporal relation between the past frames (memory) and the current frame (query), inspired by STM (Space-Time Memory) [25] for semi-supervised VOS, we use the individual convolution layer to generate the feature map (embedding), which essentially plays the role of key-value maps, so as to

reduce the model complexity. The derived feature maps reveal visual semantics for object matching and store detailed cues such as object contours for mask estimation. To determine when-and-where to retrieve related memory feature maps from, we compute similarities between the query feature and the memory features. Query feature is learned to store appearance information for decoding object mask, while memory features are learned to embed visual semantics for object matching. However, densely matching the feature maps of the query and the memory frames, requires expensive computational overheads, i.e., square time complexity *w.r.t.* the number of pixels. This motivates us to model channel-wise correlation rather than pixel-wise one, and the cost is greatly reduced by using smaller channel number. However, channel-wise attention may break down the pixel-wise similarity distribution (e.g., probability histogram), since all the pixels of pair-wise memory feature maps are taken into account channel by channel rather than pixel by pixel. Hence, we propose to make an efficient projection on the pixel-wise feature embedding of the past frames and the current frame using the similar projection skill in [11].

As illustrated in Fig. 3, the input of LTM block is the appearance feature map set $\mathcal{F}_t = \{\mathbf{F}_t^{(1)}, \mathbf{F}_t^{(2)}, \dots, \mathbf{F}_t^{(N)}, \mathbf{F}_t\} \in \mathbb{R}^{h \times w \times c_0}$ of $N$ past frames and the current frame $\mathbf{I}_t$. The feature maps of past frames and that of current frame are respectively fed into two 2D convolution layers $\phi(\cdot)$ and $\psi(\cdot)$ with $1 \times 1$ kernel, followed by reshaping the height and the width dimensions, i.e., $h \times w \to hw = n$. This results in memory features $\{\tilde{\mathbf{F}}_t^{(s)} \in \mathbb{R}^{hw \times c}\}_{s=1}^N$ ($c$ is 128) and query feature $\mathbf{Q}_t \in \mathbb{R}^{hw \times c}$, where memory features are concatenated along row dimension into one matrix, called feature memory $\mathbf{M}_t$, i.e.,

$$\mathbf{M}_t = [\tilde{\mathbf{F}}_t^{(1)}, \dots, \tilde{\mathbf{F}}_t^{(N)}] \in \mathbb{R}^{Nhw \times c}, \tag{2}$$

where $[\cdot, \cdot]$ denotes the concatenation, $N \times hw = Nhw$.

To preserve the pixel-wise similarity distribution, LTM conducts random projection on feature memory $\mathbf{M}_t$ and query feature $\mathbf{Q}_t$ at pixel level, leading to projected pixel values, i.e.,

$$m' = \frac{1}{\sqrt{\lfloor c/2 \rfloor}} \exp(\mathbf{u}^T \mathbf{m} - \frac{\|\mathbf{m}\|_2^2}{2}), \tag{3}$$

$$q' = \frac{1}{\sqrt{\lfloor c/2 \rfloor}} \exp(\mathbf{u}^T \mathbf{q} - \frac{\|\mathbf{q}\|_2^2}{2}), \tag{4}$$

where the pixel feature vector $\mathbf{m} \in \mathbb{R}^c$ is stacked in each row of feature memory matrix $\mathbf{M}_t$, and the pixel feature vector $\mathbf{q} \in \mathbb{R}^c$ is stacked in each row of query feature matrix $\mathbf{Q}_t$; the vector $\mathbf{u} \in \mathbb{R}^c$ is an orthogonal projection vector, which is randomly initialized for each projection; the constant $\lfloor c/2 \rfloor$ is a scaling factor and $\lfloor \cdot \rfloor$ rounds down fractions. Thus, we can obtain the projected pixel feature vectors $\mathbf{m}' \in \mathbb{R}^{\lfloor c/2 \rfloor}$ and $\mathbf{q}' \in \mathbb{R}^{\lfloor c/2 \rfloor}$ for each pixel in the memory feature map and the query feature map, respectively, by doing orthogonal random projections for $\lfloor c/2 \rfloor$ times as in (3). All projected pixel feature vectors are collected together to be reshaped into the matrix with the same size of that of unprojected feature matrix, i.e., $\mathbf{M}_t' \in \mathbb{R}^{Nhw \times \lfloor c/2 \rfloor}$ and $\mathbf{Q}_t' \in \mathbb{R}^{hw \times \lfloor c/2 \rfloor}$.

Therefore, the pixel-wise appearance similarity between the past frames and the query frame can be revealed by the product of projected pixel feature matrices $\mathbf{M}_t'$ and $\mathbf{Q}_t'$, i.e., $\mathbf{A}_t = \mathbf{Q}_t' \mathbf{M}_t'^{\top} \in \mathbb{R}^{hw \times Nhw}$, where the large values taken by the entries of appearance similarity matrix $\mathbf{A}_t$ indicate that the current frame shares higher similarity with the past frames in appearance. Meanwhile, its elements can be treated as attention weights of memory frames. Thus, we obtain the global feature representation by $\mathbf{G}_t = \mathbf{A}_t \mathbf{M}_t \in \mathbb{R}^{hw \times c}$, which provides the guidance to retrieve the relevant memory frames with highly similar appearance, to attend on the query frame.

Usually, the same object constantly present in video will share common appearance across frames, so learning global feature representation contributes to locating primary objects. However, the above process requires high time complexity, i.e., $\mathcal{O}(n^2 c)$, for obtaining

$$\mathbf{G}_t = \mathbf{Q}_t' \mathbf{M}_t'^{\top} \mathbf{M}_t \in \mathbb{R}^{hw \times c}. \tag{5}$$

To make the model more efficient, we first compute the channel-wise memory similarity matrix, i.e.,

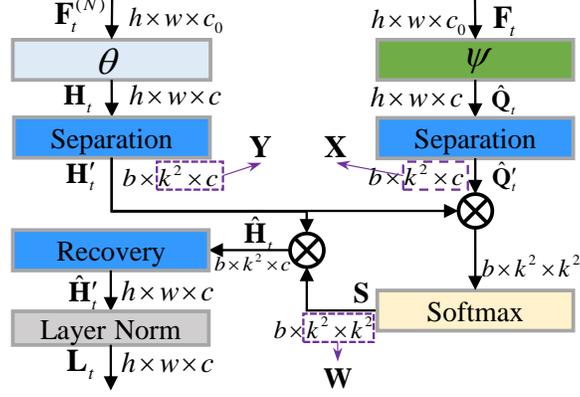$$\hat{\mathbf{A}}_t = \mathbf{M}_t'^{\top} \mathbf{M}_t \in \mathbb{R}^{\lfloor c/2 \rfloor \times c}, \tag{6}$$

Figure 4: The Short Temporal Attention (STA) block in LSTA. This block mainly consists of the 2D convolution layer with $1 \times 1$ kernel, separation and recovery operations.

where $\mathbf{M}_t \in \mathbb{R}^{Nhw \times c}$ and $\mathbf{M}'_t \in \mathbb{R}^{Nhw \times \lfloor c/2 \rfloor}$, and models the channel-wise correlation of memory features. Then, we calculate the cheap global feature representation by substituting Eq. (6) into Eq. (5), resulting in

$$\mathbf{G}_t = \mathbf{Q}'_t \hat{\mathbf{A}}_t \in \mathbb{R}^{hw \times c}, \tag{7}$$

which preserves the pixel-wise similarity distribution of the feature memory and the query feature. Especially, this formulation only requires $\mathcal{O}(nc^2)$, which can be further simplified to $\mathcal{O}(n)$ when $c \ll n$. In another word, the time complexity is linear with the number of feature map pixels, allowing the model to run very efficiently.

Besides, to avoid unstable numerical solutions due to large values, we do normalization on $\mathbf{G}_t$, i.e.,

$$\mathbf{G}'_t = \mathbf{G}_t \oslash \mathbf{D}_t \in \mathbb{R}^{hw \times c}, \tag{8}$$

where the symbol $\oslash$ denotes element-wise division (Hadamard division), and the matrix $\mathbf{D}_t = \mathbf{Q}'_t \cdot (\mathbf{M}'^{\top}_t \mathbf{1}) \in \mathbb{R}^{hw \times c}$ is used for normalization. Here, $\mathbf{1}$ is an all-one matrix with the size of $Nhw \times c$. In the end, we reshape the matrix $\mathbf{G}'_t$ to $\hat{\mathbf{G}}_t \in \mathbb{R}^{h \times w \times c}$, i.e., the global feature map, since it reflects the long-range inherent semantic relations of memory frames and the query frame.

### 3.4. Short Temporal Attention (STA)

To identify moving objects, STA block as illustrated in Fig. 4, encodes the partial-frame pixel spatiotemporal dependency of the nearest past frame and the current frame, by discovering the object motion pattern in terms of local patches of the appearance feature map. This not only helps to capture moving objects in short-term frame context, but also makes the model robust to occlusions in video.

Motivated by the attention mechanism [36] and the fact that only partial pixels in neighboring frames will change, we propose a patch-based technique called STA, for encoding temporal attention locally. Unlike the vanilla attention modeling global relations of all pixel pairs in full frame, STA models local relations of limited pixel pairs in partial frame sequentially. This is achieved by adopting the locality-based sliding window strategy, which separates one frame into a number of much smaller regions called patches. Then, STA models the local spatiotemporal relation of pixel patches between the current frame and the nearest past frame.

As illustrated in Fig. 4, the inputs of STA are the appearance feature maps of the current frame $\mathbf{I}_t$ and the nearest past frame $\mathbf{I}_t^{(N)}$, i.e., $\mathbf{F}_t$ and $\mathbf{F}_t^{(N)}$. At first, the channel dimension of $\mathbf{F}_t^{(N)}$ is reduced to $c$ by a $1 \times 1$ convolution $\theta(\cdot)$, leading to the neighbor feature map $\mathbf{H}_t = \theta(\mathbf{F}_t^{(N)}) \in \mathbb{R}^{h \times w \times c}$. For $\mathbf{F}_t$, we directly utilize its feature matrix $\mathbf{Q}_t \in \mathbb{R}^{hw \times c}$ from the convolution layer $\psi(\cdot)$ in LTM, and reshape it to query feature map $\hat{\mathbf{Q}}_t \in \mathbb{R}^{h \times w \times c}$.

STA adopts the locality-based sliding window strategy, which makes the model cheap to learn. Assume the patch size is $k$ and each feature map is separated into $b$ patches, as indicated by Fig. 5, we have $n = hw = h \times w = k \times k \times b$ pixels, and the time complexity is $\mathcal{O}(nk^2c) = \mathcal{O}(nc) = \mathcal{O}(n)$ ($k^2 \ll n$, $c \ll n$). Here we neglect the stride factor, as it does not change time complexity compared to the number of
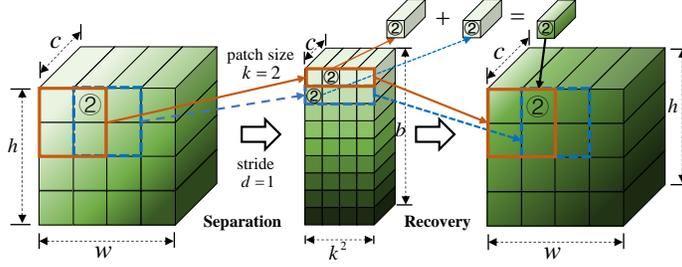
Figure 5: The illustration of separation and recovery layers in STA block.

pixels. STA models the spatial correlation of each patch with $k \times k \times c$ pixels in a feature map, and the stride $1 \leq d < k$ affects the number of patches, i.e., $b = (h - k + 1)/d \times (w - k + 1)/d$ with zero padding. Here, $k = 8$ and $d = 4$. As a result, we can obtain query patch feature tensor $\hat{\mathbf{Q}}'_t \in \mathbb{R}^{b \times k^2 \times c}$ and neighbor patch feature tensor $\mathbf{H}'_t \in \mathbb{R}^{b \times k^2 \times c}$, which are composed of $b$ patch matrices, i.e., $\{\mathbf{X}_1, \mathbf{X}_2, \ldots, \mathbf{X}_b\} \in \mathbb{R}^{k^2 \times c}$ and $\{\mathbf{Y}_1, \mathbf{Y}_2, \ldots, \mathbf{Y}_b\} \in \mathbb{R}^{k^2 \times c}$, acting as tensor slices.

To discover the pixel moving pattern in the local region of frame, STA computes the semantic similarity of query-neighbor patch pair, i.e., $(\mathbf{x}_i, \mathbf{y}_j)$, where $\mathbf{x}_i \in \mathbb{R}^c$ is stacked in the $i$-th row of $\mathbf{X}$ and $\mathbf{y}_j \in \mathbb{R}^c$ is stacked in the $j$-th row of $\mathbf{Y}$. Then, we obtain similarity value of each pixel pair by

$$\mathbf{w}_{ij} = \frac{\exp(\mathbf{x}_i^\top \mathbf{y}_j / \sqrt{c})}{\sum_j^{k^2} \exp(\mathbf{x}_i^\top \mathbf{y}_j / \sqrt{c})}. \tag{9}$$

In this way, we can compute the similarity patch by patch, and thus get the semantic similarity tensor $\mathbf{S} \in \mathbb{R}^{b \times k^2 \times k^2}$, consisting of $b$ matrices $\{\mathbf{W}_1, \mathbf{W}_2, \ldots, \mathbf{W}_b\} \in \mathbb{R}^{k^2 \times k^2}$. They help to capture the local semantic coherence of feature pairs. For those moving objects, their semantic similarity can be well encoded in the short-term spatiotemporal context, by revealing the hidden pattern. Actually, the similarity matrices play an important role in retrieving dynamics of moving objects from neighbor frame at patch level, i.e., $\{\mathbf{W}_1 \mathbf{Y}_1, \ldots, \mathbf{W}_b \mathbf{Y}_b\} \in \mathbb{R}^{k^2 \times c}$, which further act as slices of the tensor $\hat{\mathbf{H}}_t \in \mathbb{R}^{b \times k^2 \times c}$, i.e., the local feature representation.

To preserve the spatial pixel correlations, we reshape $\hat{\mathbf{H}}_t$ to the feature map $\hat{\mathbf{H}}'_t \in \mathbb{R}^{h \times w \times c}$ via the recovery layer, which is essentially the inverse process of patch separation. For each slice $k^2 \times c$, the $k^2$ entries in every column are reshaped to a local patch with the size of $k \times k$, in which way the entries in one column of all $b$ slices are reshaped to a matrix with the size of $h \times w$. Then, adding the channel dimension $c$ leads to the recovered feature map $h \times w \times c$. Due to the stride of sliding window, there is redundant samplings on feature map. Hence, we simply sum the elements of those local patch features which are projected onto the same pixel, resulting in the feature map with the same size of the original one. However, those large sum values might possibly lead to numerical instability of the data distribution. To deal with this issue, we impose the layer normalization on $\hat{\mathbf{H}}'_t$, i.e., normalization across all channels of feature map. As a result, we obtain the local feature map $\mathbf{L}_t \in \mathbb{R}^{h \times w \times c}$, which encodes the short-range spatiotemporal pixel relations of the nearest neighbor frame and the query frame in terms of small patches.

### 3.5. Decoder

Decoder is composed of convolution layer, anisotropic convolution (AIC) [14] block (consisting of several 2D convolution layers), and up-sampling operation. Its goal is to consider both the long-range and the short-range spatiotemporal pixel correlations of the past frames and the current frame, by make a fusion on three data streams, i.e., global feature map $\hat{\mathbf{G}}_t$, local feature map $\mathbf{L}_t$, and query feature map $\hat{\mathbf{Q}}_t$.

Particularly, we first concatenate the above three feature maps along the channel dimension into a unified feature map with size of $h \times w \times 3c$, whose channel dimension is reduced to $c$ via one $3 \times 3$ 2D convolution layer. Then, the unified feature map is fed into the AIC module that helps to discriminate objects from the spatial context. This is followed by passing the other $3 \times 3$ 2D convolution layer for reducing the channel dimension into 2. Hereafter, the resolution of feature map is enlarged to that of original RGB frame by

9

bilinear up-sampling. Finally, it yields the object-like pixel probability using softmax function, namely

$$\tilde{\mathbf{P}}_t = \Theta(\hat{\mathbf{G}}_t, \mathbf{L}_t, \hat{\mathbf{Q}}_t) \in \mathbb{R}^{H \times W \times 2}, \tag{10}$$

where $\Theta(\cdot)$ denotes Decoder. The elements of the first (index 0) and the second (index 1) channel denote the probabilities of pixels belonging to background and primary object, respectively. This probability tensor can be easily transformed to a binary matrix by taking the channel index of the higher probability for each pixel, i.e., $\hat{\mathbf{P}}_t = \{0, 1\} \in \mathbb{R}^{H \times W}$.

### 3.6. Loss Function

To optimize the LSTA model, we adopt the Cross-Entropy (CE) loss with the online hard example mining strategy [32]. It selects those hard pixels (usually those with larger loss values) to calculate the CE loss, which is beneficial for promoting the robustness of discriminating those ambiguous pixel regions. The model loss is computed by

$$\mathcal{L}_1(\tilde{\mathbf{P}}_t, \mathbf{P}_t) = \overline{\mathrm{Max}}_r(\{-p_{tz}^0 \log \tilde{p}_{tz}^0 - p_{tz}^1 \log \tilde{p}_{tz}^1\}_{z=1}^{HW}), \tag{11}$$

where $\overline{\mathrm{Max}}_r(\cdot)$ ($r = \lfloor \frac{HW}{16} \rfloor$ and $HW = H \times W$) means taking the average of those $r$ largest loss values of pixels, $\{\tilde{p}_{tz}^0, \tilde{p}_{tz}^1\}$ are the predicted probability values of the $z$-th pixel in current frame $\mathbf{I}_t$, and $\{p_{tz}^0, p_{tz}^1\}$ are the corresponding ground-truth object mask values.

Inspired by knowledge distillation, we utilize the semi-supervised VOS model, i.e., STM [25], as teacher network, which provides guidance for our unsupervised VOS model, i.e., LSTA, as student network. Note that, the STM model trained on DAVIS17 [27] and YouTube-VOS [45] does not participate in training our LSTA model but only helps to yield the initial pixel probability of frames, i.e., soft labels, to guide the loss computation. And the inference process does not involve teacher network as well. Assume that the initial pixel probability from teacher network is $\bar{\mathbf{P}} \in \mathbb{R}^{H \times W \times 2}$, we compute the following loss:

$$\mathcal{L}_2(\tilde{\mathbf{P}}_t, \bar{\mathbf{P}}_t) = \frac{1}{HW} \sum_{z=1}^{HW} (-\bar{p}_{tz}^0 \log \tilde{p}_{tz}^0 - \bar{p}_{tz}^1 \log \tilde{p}_{tz}^1), \tag{12}$$

where $\{\bar{p}_{tz}^0, \bar{p}_{tz}^1\}$ are the initial probability values of the $z$-th pixel in the current frame $\mathbf{I}_t$.

Therefore, the total loss of our LSTA model is

$$\mathcal{L}(\tilde{\mathbf{P}}_t, \mathbf{P}_t, \bar{\mathbf{P}}_t) = \alpha \mathcal{L}_1(\tilde{\mathbf{P}}_t, \mathbf{P}_t) + (1 - \alpha) \mathcal{L}_2(\tilde{\mathbf{P}}_t, \bar{\mathbf{P}}_t), \tag{13}$$

where the tradeoff parameter $\alpha > 0$ is used to balance the contribution of the two loss terms to the objective function. Here, we use an empirical value 0.5.

To summarize our approach, we briefly show the training process in Algorithm 1 and the inference process in Algorithm 2. Our method is relevant to the widely used LSTM and Transformer, which captures the short-term and the long-term temporal relations, respectively. However, their working mechanisms are different. In particular, LSTM adopts the gating skill to control the history information in a sequence, while our STA block uses the attention mechanism to discover the object motion pattern in terms of local patches of the appearance feature map, which encodes the partial-frame pixel spatiotemporal dependency with the locality-based sliding window strategy. Moreover, Transformer adopts the self-attention to model the global dependency and needs the square time complexity, while our LTM block encodes the full-frame pixel spatiotemporal dependency in terms of appearance similarity at the cost of almost linear complexity by computing the channel-wise memory similarity matrix.

## 4. Experiments

This section shows comprehensive empirical studies on several benchmark data sets. All experiments were conducted on a machine with NVIDIA Titan RTX Graphics Card for training and NVIDIA Titan Xp for inference, and our LSTA model was compiled using PyTorch 1.8, Python 3.6, and CUDA 11.0.

---
**Algorithm 1:** Training Process of LSTA model.
---
**Input:** Training videos; ground truth mask $\mathcal{P}$; model parameters $\Omega$; number of past frames $N$;
$\quad\quad\quad \alpha = 0.5; Iter_{max} = 7.5e4$.

1  Randomly select $T$ frames as query for each video.
2  **while** *not converged* **do**
3  $\quad$ Randomly select one query frame $\mathbf{I}$ from one video in sequence.
4  $\quad$ Select $N$ past frames $\mathbf{I}^{(1)}, \mathbf{I}^{(2)}, ..., \mathbf{I}^{(N)}$ .
5  $\quad$ Feed query frame and past frames into Encoder to obtain appearance feature map $\mathbf{F}$ and
$\quad\quad \{\mathbf{F}^{(1)}, \mathbf{F}^{(2)}, \ldots, \mathbf{F}^{(N)}\}$.
6  $\quad$ Input appearance feature maps $\{\mathbf{F}, \mathbf{F}^{(1)}, \mathbf{F}^{(2)}, ..., \mathbf{F}^{(N)}\}$ into LTM to obtain global feature map
$\quad\quad \hat{\mathbf{G}}$.
7  $\quad$ Feed appearance feature maps $\mathbf{F}$ and $\mathbf{F}^{(N)}$ into STA to obtain local feature map $\mathbf{L}$ and query
$\quad\quad$ feature map $\hat{\mathbf{Q}}$.
8  $\quad$ Obtain the object-like pixel probability $\tilde{\mathbf{P}}$ using $\hat{\mathbf{G}}$, $\mathbf{L}$ and $\hat{\mathbf{Q}}_t$ as in Eq. (10).
9  $\quad$ Calculate the first loss $\mathcal{L}_1$ as in Eq. (11) and the second loss $\mathcal{L}_2$ as in Eq. (12).
10 $\quad$ Calculate the total loss function $\mathcal{L}$ in Eq. (13).
11 $\quad$ Update model parameters $\Omega$ using SGD.
12 **end**
**Output:** Trained model.
---

---
**Algorithm 2:** Inference Process of LSTA model.
---
**Input:** Video frames; trained model.
1  Obtain appearance feature map for every video frame by Encoder.
2  Feed appearance feature maps of query frame and past frames to LTM to obtain global feature map.
3  Feed appearance feature maps of query frame and nearest neighbor frame to STA to obtain local
$\quad$ feature map and query feature map.
4  Input global feature map, local feature map, and query feature map to Decoder, which yields object
$\quad$ mask.
**Output:** Object mask set.
---

## 4.1. Data Sets

In total, there are four publicly available VOS data sets used in the experiments. Details are shown below.

**DAVIS**[1] short for Densely Annotated Video Segmentation, provides two kinds of frame resolution, i.e., 480p and 1080p, with pixel-level frame mask. It has two versions, i.e., DAVIS2016[26] and DAVIS2017[27], involving various scenes, such as animal, sports, and traffic vehicles. The former has 50 video sequences, which are divided into 30 training videos and 20 validation videos for inference; each video contains only one object, and there are 3,455 frames with ground-truth masks. The latter DAVIS2017 is an expansion of the former, and the number of videos increases to 150, among which there are 90 videos (60 for training and 30 for validation) with 10,459 frames with ground-truth masks; each video may contain more than one objects, adding difficulty to the task, and there are totally 376 objects.

**YouTube-VOS**[2] [28] collects video clips from YouTube web site, including various classes, such as animal, transportation, accessory, and human event. Each clip usually contains multiple objects, with a duration of 3s to 6s. It has three subsets, and we only use its training set, including 3,471 videos with dense (6 fps) object annotations, 65 categories, and 5,945 unique object instances.

**YouTube-Objects**[3] [28] is composed of 126 videos collected from YouTube by querying for the names of 10 object classes. The duration of each video varies between 30s and 180s, and each video contains one

---

[1] https://davischallenge.org/index.html
[2] https://youtube-vos.org/dataset/vos/
[3] https://data.vision.ee.ethz.ch/cvl/youtube-objects/

object. Following [55][40], we use all videos for inference.

**FBMS**[4] [24] short for Freiburg-Berkeley Motion Segmentation, contains 59 video sequences, which are separated to 29 training videos and 30 validation videos, each of which has one object. There are 720 frames with pixel-level mask annotations, which are made with an interval of 20 frames. Following [55][20], we only use validation set for inference.

### 4.2. Evaluation Metrics

We evaluate our LSTA model on DAVIS, YouTube-Objects, and FBMS benchmarks, and the model is learned using the training sets of DAVIS2017 and YouTube-VOS. Following the previous works [55][20][48], we use region similarity $\mathcal{J}$, contour accuracy $\mathcal{F}$, and $\mathcal{J}\&\mathcal{F}$ as the evaluation criteria. For DAVIS, we used the official benchmark code [26]. For YouTube-Objects and FBMS, we use $\mathcal{J}$ Mean as the metric.

Region similarity $\mathcal{J}$ is the Intersection over Union (IoU, namely Jaccard coefficient) of predicted mask $\hat{\mathbf{P}}$ and ground-truth mask $\mathbf{P}$, which reflects the spatial mask accuracy and is a frame size-agnostic metric. It is computed by $\mathcal{J} = \frac{|\hat{\mathbf{P}}_t \cap \mathbf{P}_t|}{|\hat{\mathbf{P}}_t \cup \mathbf{P}_t|}$. Contour accuracy $\mathcal{F}$ estimates whether the contour of predicted mask $\hat{\mathbf{P}}$ is similar with that of ground-truth mask $\mathbf{P}$. From a contour perspective, one can interpret $\hat{\mathbf{P}}$ and $\mathbf{P}$ as a set of closed contours $\mathcal{C}(\hat{\mathbf{P}})$ and $\mathcal{C}(\mathbf{P})$ delimiting the spatial extent of the mask. So one can compute the contour-based precision $P_c$ and recall $R_c$ via a bipartite graph matching. The $\mathcal{F}$ measure is a harmonic value, i.e., $\mathcal{F} = \frac{2P_c R_c}{P_c + R_c}$.

$\bar{\mathcal{J}}$ denotes $\mathcal{J}$ Mean and $\bar{\mathcal{F}}$ denotes $\mathcal{F}$ Mean, each of which is the average result over all test videos. Meanwhile, we use the mean value of region similarity and contour accuracy as the overall evaluation metric, i.e., $\overline{\mathcal{J}\&\mathcal{F}}$ ($\mathcal{J}\&\mathcal{F}$ Mean) over all videos. In addition, we use Frame Per Second (FPS) as the metric to evaluate the inference speed.

### 4.3. Experimental Setup

**Training Phase**. The Encoder of LSTA is initialized by the DeepLab v3+ [3] model pre-trained on MS COCO, while the other modules are randomly initialized using Xavier. In each iteration, we randomly sample a single frame from each of 4 videos (batch size is 4) as query frame, and its all previous frames are grouped into $N = 5$ bins, from each of which we randomly select one frame, resulting in $N$ past frames with temporal relations. When going through all available training videos once, it finishes one epoch. For each frame, it is randomly cropped to $465 \times 465 \times 3$, while random horizontal flipping and scaling are applied. The maximum iteration number is $7.5e4$, and we adopt the SGD (Stochastic Gradient Descent) optimizer with a momentum of 0.9, a weight decay of $1.5e$-4, and an initial learning rate of $6e$-3. Note that our model is trained on DAVIS2017[27] and YouTube-VOS[45] with $5e4$ iterations, and is fine-tuned on DAVIS2017[27] with $2.5e4$ iterations to further boost the generalization ability.

**Inference Phase**. For unseen videos, according to Algorithm 2, LSTA sequentially takes the current frame as query frame and previous $N$ frames as past frames without any object prior, and outputs the corresponding object masks. Note that there will be insufficient past frames for the foremost $N$ query frames. For such cases, we use the succeeding frames to compensate for the lacking past frames. In addition, we follow AD-Net [48] to use instance pruning to filter out some cluttered background by employing instance bounding boxes.

### 4.4. Quantitative Results

We show the quantitative comparison results on DAVIS2016, DAVIS2017, YouTube-Objects, and FBMS, with rigorous analysis in the following.

**DAVIS2016**. The results of our LSTA model and a number of state-of-the-art alternatives are recorded in Table 1. Among them, the above seven methods are semi-supervised models, including FEELVOS (Fast End-to-End Embedding Learning) [38], AGUNet (Annotation-Guided U-Net)[49], MVOS-OL (Meta VOS Online Learning)[44], and DDEAL (Directional Deep Embedding and Appearance Learning) [50]. The remaining ones are all unsupervised models, including PDB (Pyramid Dilated Bidirectional ConvLSTM)[33], AGS (Attention-Guided object Segmentation)[41], AGNN (Attentive Graph Neural Networks)[40], COSNet

---

[4]https://lmb.informatik.uni-freiburg.de/resources/datasets/moseg.en.html

Table 1: Performance comparisons on DAVIS2016. The methods in top group are semi-supervised and the rest are unsupervised. att: attention mechanism; flow: optical flow feature; pp: post-processing; crf: conditional random field skill; ip: instance pruning; '-' means the speed is unavailable; $\dagger$ denotes the fps on RTX 2080Ti GPU.

| Method | Venue | att | flow | pp | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ | FPS |
|---|---|---|---|---|---|---|---|---|
| AGUNet[49] | PR'21 | ✓ | | | 80.7 | 81.0 | 80.9 | 11.1 |
| FEELVOS[38] | CVPR'19 | | | | 81.1 | 82.2 | 81.7 | 2.0 |
| MVOS-OL[44] | TPAMI'20 | | | | 83.3 | 84.1 | 83.7 | 2.3 |
| DDEAL[50] | TNNLS'21 | | | | 85.1 | 85.7 | 85.4 | 25.0 |
| PDB[33] | ECCV'18 | | | crf | 77.2 | 74.5 | 75.9 | - |
| AGNN[40] | ICCV'19 | ✓ | | crf | 80.7 | 79.1 | 79.9 | 0.3 |
| COSNet[20] | CVPR'19 | ✓ | | crf | 80.5 | 79.5 | 80.0 | 1.2 |
| STEm-Seg[1] | ECCV'20 | | | | 80.6 | 80.6 | 80.6 | 0.7 |
| AD-Net[48] | ICCV'19 | ✓ | | ip | 81.7 | 80.5 | 81.1 | 0.3 |
| MATNet[55] | TIP'20 | ✓ | ✓ | crf | 82.4 | 80.7 | 81.6 | 1.3 |
| FSNet[9] | ICCV'21 | ✓ | ✓ | | 82.3 | 83.3 | 82.8 | 12.5 |
| FSNet*[9] | ICCV'21 | ✓ | ✓ | crf | 83.4 | 83.1 | 83.3 | - |
| 3DC[22] | BMVC'20 | | | | 83.9 | 84.2 | 84.1 | 4.5 |
| DASPP[53] | PR'21 | | | | 63.4 | 60.2 | 61.8 | 29.4 |
| AGS[41] | TPAMI'21 | ✓ | | crf | 79.7 | 77.4 | 78.6 | - |
| RTNet[30] | CVPR'21 | ✓ | ✓ | | 85.6 | 84.7 | 85.2 | 4.3 |
| OFS[23] | TPAMI'23 | | ✓ | | 69.3 | 70.7 | 70.0 | - |
| FEM-Net[57] | TCSVT'22 | ✓ | ✓ | | 79.9 | 76.9 | 78.4 | 16.0 |
| IMCNet[43] | TCSVT'22 | ✓ | | | 82.7 | 81.1 | 81.9 | - |
| IMP[13] | AAAI'22 | ✓ | | | 84.5 | **86.7** | **85.6** | 1.79$^\dagger$ |
| TMO[6] | WACV'23 | | ✓ | | **85.6** | 86.6 | 86.1 | 24.8 |
| LSTA (Ours) | | ✓ | | | 82.4 | 84.3 | 83.4 | **42.8** |
| LSTA*(Ours) | | ✓ | | ip | 82.7 | 84.8 | 83.8 | 36.2 |

Table 2: Computational analysis on DAVIS2016.

| Method | Venue | Backbone | $\overline{\mathcal{J}\&\mathcal{F}}$ | Params(M)↓ | FLOPs(G)↓ | Speed (FPS) | |
|---|---|---|---|---|---|---|---|
| | | | | | | TITAN Xp↑ | 2080Ti↑ |
| MATNet[55] | TIP'20 | ResNet101 | 81.6 | 142.7 | **193.7** | 1.3 | 1.8 |
| COSNet[20] | CVPR'19 | DeepLabv3 | 80.0 | 81.2 | 585.5 | 1.2 | 1.6 |
| RTNet[30] | CVPR'21 | ResNet101 | 85.2 | 277.2 | 489.6 | 4.3 | 5.8 |
| IMCNet[43] | TCSVT'22 | ResNet101 | 81.9 | **47.9** | 391.8 | 11.0 | 19.0 |
| TMO[6] | WACV'23 | ResNet101 | **86.1** | 93.0 | 344.9 | 24.8 | 39.1 |
| LSTA(Ours) | - | DeepLabv3 | 83.4 | 60.5 | 349.5 | **42.8** | **62.3** |

Table 3: Component computational analysis on DAVIS2016.

| Blocks | Encoder | LTM | STA | Decoder | LSTA |
|---|---|---|---|---|---|
| FLOPs(G) | 279.5 | 4.2 | 2.4 | 63.4 | 349.5 |
| Params(M) | 59.2 | 0.03 | 0.02 | 1.2 | 60.5 |

(Co-Attention Siamese Networks)[20], STEm-Seg (Spatio-Temporal Embeddings for instance Segmentation) [1], AD-Net (Anchor Diffusion Network)[48], MATNet (Motion-Attentive Transition Network)[55], FSNet (Full-duplex Strategy Network)[9], FEM-Net (Flow Edge-based Motion-attentive Network)[57], 3DC (3D Convolutions)[22], RTNet (Reciprocal Transformation Network)[30], IMCNet (Implicit Motion-Compensated Network) [43], OFS (Optical Flow Segmentation), [23], DASPP (Dynamic Astrous Spatial Pyramid Pooling) [53], IMP (Iterative Mask Propagation)[13], and TMO (Treating Motion as Option)[6]. Among them, many methods such as OFS, FEM-Net, RT-Net and TMO, adopt the two-stream VOS framework that employs op-

tical flow as the motion modality to capture the temporal relations; IMCNet aligns motion information from nearby frames to the current frame, and adopts a co-attention mechanism to learn the global representation; IMP repeats easy frame selection with mask propagation but the inference speed is very low.

These records show that our LSTA model enjoys more satisfying overall performance, and very competitive with SOTA unsupervised alternatives in terms of $\mathcal{F}$ Mean, i.e., $84.3\%$. Especially, LSTA ranks Top-1 in the inference speed, achieving 42.8 fps, almost 9.5 times of that of the best candidate 3DC. This demonstrates that LSTA can be deployed in highly demanding applications with its real-time segmentation ability. Besides, we have the following observations:

- Our LSTA can beat against several semi-supervised VOS methods, such as FAVOS and RGMP, in terms of both region similarity and contour accuracy. Meanwhile, LSTA is much faster ($1.7\times$) than DDEAL who has the best performance in semi-supervised setting.

- Most of the unsupervised VOS methods adopt post-processing techniques, such as conditional random field or instance pruning to refine the object mask. We also use instance pruning to slightly upgrade the performance by $0.4\%$ of $\mathcal{J}\&\mathcal{F}$ Mean.

- Unlike existing unsupervised VOS methods that use costly optical flow features, our model attempts to discover the patterns of constantly present and moving objects, by encoding both global and local spatiotemporal correlations across frames with the long-short temporal attention mechanism.

- While the latest work TMO achieves the best segmentation performance, it requires the expensive optical flow features and the inference speed is much lower than ours, i.e., 24.8fps versus 42.8fps. This demonstrates the LSTA strikes a good balance between the performance and speed.

In addition, we provide the computational comparison and component computational analysis in Table 2 and Table 3, respectively. Note that we also provide the inference speed using 2080Ti card for reference. It can be seen from the tables that our LSTA achieves the best inference speed at 42.8fps using TITAN Xp, which is almost the twice faster than that of the second best one. Meanwhile, the developed LTM module and STA module are very lightweight in terms of FLOPs with the negligible parameter size.

Table 4: Performance comparisons on DAVIS2017.

| Method | Venue | att | flow | pp | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|---|---|---|---|
| RVOS[37] | CVPR'19 | | | | 36.8 | 45.7 | 41.2 |
| PDB[33] | ECCV'18 | | | crf | 53.2 | 57.0 | 55.1 |
| MATNet[55] | TIP'20 | ✓ | ✓ | crf | 56.7 | 60.4 | 58.6 |
| STEm-Seg[1] | ECCV'20 | | | | 61.5 | 67.8 | 64.7 |
| UnOVOST[21] | WACV'20 | | ✓ | | 66.4 | 69.3 | 67.9 |
| AGS[41] | TPAMI'21 | ✓ | | crf | 55.5 | 59.5 | 57.5 |
| DyStaB[47] | CVPR'21 | | ✓ | crf | 58.9 | - | 58.9 |
| TAODA[54] | CVPR'21 | | | | 63.7 | 66.2 | 65.0 |
| LSTA (Ours) | | ✓ | | | **70.8** | **75.8** | **73.3** |
| LSTA$^{\dagger}$ (Ours) | | ✓ | | | <u>67.8</u> | <u>72.3</u> | <u>70.1</u> |

**DAVIS2017**. This dataset is much more difficult for segmentation, since there may exist multiple objects in a single video. To handle this case, some previous unsupervised VOS works, such as UnOVOST (Unsupervised Offline VOS) [21], MATNet (Motion-Attentive Transition Network) [55], AGS (Attention-Guided object Segmentation) [41], employ instance segmentation model Mask R-CNN [8] to obtain object proposals involving mask and boundary box. Unlike them, TAODA (Target-Aware Object Discovery and Association) [54] introduces an instance discrimination network to obtain object proposals in a bottom-up fashion. Moreover, DyStaB (Dynamic-Static Bootstrapping) [47] employs a motion segmentation module to perform temporally consistent region separation, and it requires the expensive optical flow features and CRF post-processing to get final results. Note that STEm-Seg [1] used spatiotemporal embeddings to find instances, which is based on Gaussian distribution, failing to handle complex appearance of objects. The

earlier works, PDB (Pyramid Dilated Bidirectional ConvLSTM) [33] and RVOS (Recurrent network) [37] identify the instance by recurrent neural networks with ConvLSTM, which fails to encode long-range temporal relations. Similarly, we use HTC (Hybrid Task Cascade) [2] to obtain object proposals of the first frame, which are then processed to yield object mask, and adopt STCN (Space-Time Correspondence Network) [4] to obtain initial object masks of subsequent frames in a semi-supervised manner. After that, we use our LSTA model to obtain predicted masks that contain multiple objects. In addition, we follow UnOVOST using Mask R-CNN to extract object proposals and its merging strategy with the pixel probability obtained by our model, and the results are listed in the bottom row.

The comparison results of the above methods are tabulated in Table 4, which has shown the significant improvements brought by our model with good generalization ability. For example, LSTA achieves $70.8\%$ on $\mathcal{J}$ Mean, $75.8\%$ on $\mathcal{F}$ Mean, and $73.3\%$ on $\mathcal{J}\&\mathcal{F}$ Mean, which have improvements of $4.4\%$, $6.5\%$, and $5.4\%$ compared to the most competitive alternative, i.e., UnOVOST. We attribute this to the fact that our approach is able to encode both long-range and shot-range spatiotemporal pixel-wise relations of the current frame and the past frames, which helps to better capture constantly present objects and moving objects in video.

Table 5: Performance comparisons on YouTube-Objects with 10 categories. The number of videos in each category is in parenthesis.

| Method | SegFlow[5] | PDB[33] | MATNet[55] | COSNet[20] | AGNN[40] | AGS[41] | RTNet[30] | IMCNet[43] | TMO[6] | LSTA |
|---|---|---|---|---|---|---|---|---|---|---|
| Venue | ICCV'17 | ECCV'18 | TIP'20 | CVPR'19 | ICCV'19 | TPAMI'21 | CVPR'21 | TCSVT'22 | WACV'23 | Ours |
| att | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ |
| flow | ✓ | | ✓ | | | | ✓ | | ✓ | |
| pp | | crf | crf | crf | crf | crf | | | | |
| Airplane(6) | 65.6 | 78.0 | 72.9 | 81.1 | 81.1 | **87.7** | 84.1 | 81.1 | <u>85.7</u> | 85.1 |
| Bird(6) | 65.4 | 80.0 | 77.5 | 75.7 | 75.9 | 76.7 | <u>80.2</u> | **81.1** | 80.0 | 75.9 |
| Boat(15) | 59.9 | 58.9 | 66.9 | <u>71.3</u> | 70.7 | **72.2** | 70.1 | 70.3 | 70.1 | 66.5 |
| Car(7) | 64.0 | 76.5 | <u>79.0</u> | 77.6 | 78.1 | 78.6 | **79.5** | 77.1 | 78.0 | 78.4 |
| Cat(16) | 58.9 | 63.0 | **73.7** | 66.5 | 67.9 | 69.2 | 71.8 | 73.3 | <u>73.6</u> | 72.4 |
| Cow(20) | 51.2 | 64.1 | 67.4 | 69.8 | 69.7 | 64.6 | <u>70.1</u> | 66.8 | **70.3** | 67.1 |
| Dog(27) | 54.1 | 70.1 | 75.9 | 76.8 | <u>77.4</u> | 73.3 | 71.3 | 74.8 | 76.8 | **77.9** |
| Horse(14) | 64.8 | <u>67.6</u> | 63.2 | 67.4 | 67.3 | 64.4 | 65.1 | 64.8 | 66.2 | **68.5** |
| Motorbike(10) | 52.6 | 58.4 | 62.6 | <u>67.7</u> | **68.3** | 62.1 | 64.6 | 58.7 | 58.6 | 65.5 |
| Train(5) | 34.0 | 35.3 | 51.0 | 46.8 | 47.8 | 48.2 | 53.3 | <u>56.8</u> | 47.0 | **57.5** |
| $\overline{\mathcal{J}}$ | 57.1 | 65.5 | 69.0 | 70.5 | 70.8 | 69.7 | <u>71.0</u> | 70.5 | **71.5** | **71.5** |
| FPS | - | - | - | 0.9 | 0.4 | 4.2 | - | - | <u>18.5</u> | 35.5 |

**YouTube-Objects**. Following previous work[55], we give segmentation results in terms of $\mathcal{J}$ Mean for each of the *ten* semantic categories, as in Table 5. The compared methods include SFL[5], PDB[33], MATNet[55], AGS[41], COSNet[20], AGNN[40], RTNet[30], IMCNet [43], and TMO[6]. Among them, our model outperforms the rest ones in overall performance, achieving $71.5\%$ on $\mathcal{J}$ Mean at 35.5 fps in inference speed. Especially, the performance of our method is comparable to the latest TMO model which employs the more complicated framework, and our speed is twice faster. The encouraging records have verified the advantage of our approach in terms of both effectiveness and efficiency.

Table 6: Comparison results on FBMS validation set.

| Method | OBN[16] | PDB[33] | COSNet[20] | MATNet[55] | OFS[23] | AGS[41] | DASPP[53] | LSTA |
|---|---|---|---|---|---|---|---|---|
| Venue | ECCV'18 | ECCV'18 | CVPR'19 | TIP'20 | TPAMI'23 | TPAMI'21 | PR'21 | Ours |
| att | | | ✓ | ✓ | | ✓ | ✓ | ✓ |
| flow | ✓ | | | ✓ | ✓ | | | |
| pp | crf | crf | crf | crf | | crf | | |
| $\overline{\mathcal{J}}$ | 73.9 | 74 | 75.6 | <u>76.1</u> | 57.8 | 76.0 | 62.3 | **77.3** |
| FPS | <u>4.7</u> | - | 0.9 | 0.4 | - | - | - | **41.4** |

**FBMS**. We compare our model with OBN (Object Bilateral Networks)[16], PDB[33], COSNet[20], MATNet[55], OFS [23], AGS[41], and DASPP [53] methods on FBMS, whose results are shown in Table 6. From these records, we see our LSTA method gets the highest region similarity, i.e., $77.3\%$, with a margin

Table 7: Ablation study of individual components on DAVIS2016.

| Blocks | $\mathcal{J}$ Mean | $\mathcal{F}$ Mean | $\overline{\mathcal{J}\&\mathcal{F}}$ | Gain |
|---|---|---|---|---|
| Baseline | 76.3 | 79.6 | 77.9 | - |
| +LTM | 80.1 | 82.3 | 81.2 | +3.3 |
| +STA | 79.5 | 82.5 | 81.0 | +3.1 |
| +LTM+STA | 81.4 | 83.8 | 82.6 | +4.7 |
| +LTM+STA+$\mathcal{L}_2$ | **82.4** | **84.2** | **83.3** | +5.4 |



Figure 6: Visualization results of the ablation study on individual components. Row 1/3: LTM, Row 2/4: LTM+STA, Row 3/6: LTM+STA+$\mathcal{L}_2$. The dashed yellow circle highlights the difference.

of 1.2% compared to the second best, i.e., MATNet. Meanwhile, our method achieves a satisfying balance between segmentation performance and inference speed without optical flow features and post-processing, e.g., its inference speed is 41.4fps, which helps to being deployed in highly-demanding environment.

### 4.5. Ablation Study

This section makes extensive analysis on the contribution of individual components of LSTA, the number of past frames $N$ and the strategy of selecting them, the patch size and stride in STA block, the influences of channel numbers $c$ and whether sharing Conv2D layers, the tradeoff parameter $\alpha$, and the performance on the videos with various visual characteristics on DAVIS2016.

**Individual components**. We show the effects of different components in Table 7, which examines the Light Temporal Attention block and the Short Temporal Attention block as well as $\mathcal{L}_2$ loss terms. For baseline, we use $\mathcal{L}_1$ loss term. From the table, we observe that both LTM and STA promote the performance by 3.8% (row 2) and 3.2% (row 3), respectively, on $\mathcal{J}$ Mean, compared to baseline. Besides, the coupling of LTM and STA brings about an upgrade (row 4) by some margin of 0.6% and 3.2% on $\mathcal{J}\&\mathcal{F}$ Mean, respectively. This demonstrates the necessity of simultaneously considering both global and local spatiotemporal pixel relations of the current frame and the past frames. In addition, the knowledge distillation skill is beneficial for boosting the performance by 1.0% (bottom row) on $\mathcal{J}$ Mean. Meanwhile, we show some visualization examples of using different components in Fig. 6. As shown in the figure, the error area becomes smaller or disappears when using both STA module and LTM module (Row 2 and 4); meanwhile, the segmentation quality is further improved by adding the knowledge distillation loss $\mathcal{L}_2$ (Row 3 and 6).

**Number of past frames in LTM**. We vary the number of past frames from 2 to 11 for LTM, and the results are shown in Table 8. It can be seen that when $N$ is 5, the model achieves the best performance, which

Table 8: Number $N$ of past frames.

| $N$ | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|
| 2 | 82.0 | 84.1 | 83.1 |
| 3 | 82.1 | 84.1 | 83.1 |
| 5 | **82.4** | **84.2** | **83.3** |
| 7 | 82.1 | 84.1 | 83.1 |
| 9 | 81.9 | 84.1 | 83.0 |
| 11 | 81.9 | 84.0 | 83.0 |

Table 9: Past frame selection ($N$=5).

| Frames | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|
| first | 81.2 | 83.5 | 82.4 |
| prev | 82.0 | 84.1 | 83.0 |
| first&prev | 81.4 | 83.8 | 82.6 |
| rand $N$ | 81.9 | 84.0 | 82.9 |
| every $N$ | 81.9 | 84.1 | 83.0 |
| prev $N$ | **82.4** | **84.3** | **83.4** |

Table 10: Patch size $k$ and stride $d$ in STA block.

| $k$ | $d$ | $\mathcal{J}$ Mean | $\mathcal{F}$ Mean | $\mathcal{J}\&\mathcal{F}$ Mean |
|---|---|---|---|---|
| 2 | 1 | 82.0 | **84.2** | 83.1 |
| 4 | 2 | 82.0 | **84.2** | 83.1 |
| 8 | 4 | **82.4** | **84.3** | **83.4** |
| 16 | 8 | 82.0 | 84.1 | 83.0 |
| 32 | 16 | 81.9 | 83.9 | 82.9 |

Table 11: Channel Number $c$.

| $c$ | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|
| 16 | 79.4 | 80.9 | 80.2 |
| 32 | 81.4 | 83.7 | 82.6 |
| 64 | **82.4** | **84.3** | **83.4** |
| 128 | 81.8 | 84.0 | 82.9 |
| 256 | 81.5 | 83.9 | 82.7 |

Table 12: Sharing Conv2D layer.

| Share | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|
| $\psi = \phi$ | 81.5 | 83.9 | 82.7 |
| $\psi = \theta$ | 81.6 | 84.2 | 82.9 |
| $\phi = \theta$ | 81.7 | 84.1 | 82.9 |
| $\psi = \phi = \theta$ | 81.6 | 84.1 | 82.9 |
| None | **82.4** | **84.3** | **83.4** |

suggests that much more past frames can not provide additional spatiotemporal information due to frame redundancy.

**Selecting past frames**. Table 9 gives the results of six ways of selecting the past frames. Selecting the nearest frame (row 2) performs better than using the first frame (row 1). Besides, using both the first frame and the previous one will slightly degrade performance in comparison of using only the previous one, which might be reason that the object mask changes a lot with time and the first frame may mislead mask prediction. When using more past frames, e.g., two frames in row 3 and five frames in the last three rows, the segmentation performance is improved, especially using previous $N$ nearest frames (bottom row) in inference phase.

**Patch size and stride in STA**. STA adopts the locality-based sliding window strategy to reduce the computational cost, which is largely governed by the patch size $k$ and the stride $d$. We vary them from 2 to 32 and 1 to 16, respectively, and the results are recorded in Table 10. From the table, STA block works the best when the batch size is 8 and the stride is 4, which are both modest values. Larger patches or smaller ones do not help the improvements of capturing local spatiotemporal pixel relations of the current frame and the previous frame.

**Number of channels in Conv2D**. Table 11 shows the results of increasing the number of channels ($c$) in Conv2D from 16 to 256. The region similarity metric is improved by 3.2% using 64 channels, compared to that with 16 channels. More channels encode more accurate spatial structure of frame data, but the performance degrades when $c$ is over 100. This might because some noise in channel dimension is mixed with the feature maps.

**Sharing Conv2D layer**. Our model adopts Conv2D layers for encoding the past frames in LTM by $\phi(\cdot)$, the previous frame in STA by $\theta(\cdot)$, and the current frame in STA by $\psi(\cdot)$. So we explore the influences of sharing them in various forms as in Table 12. When sharing either two of them or three all, the performance is worse than using Conv2D layers independently for them. It suggests that learning parameters independently
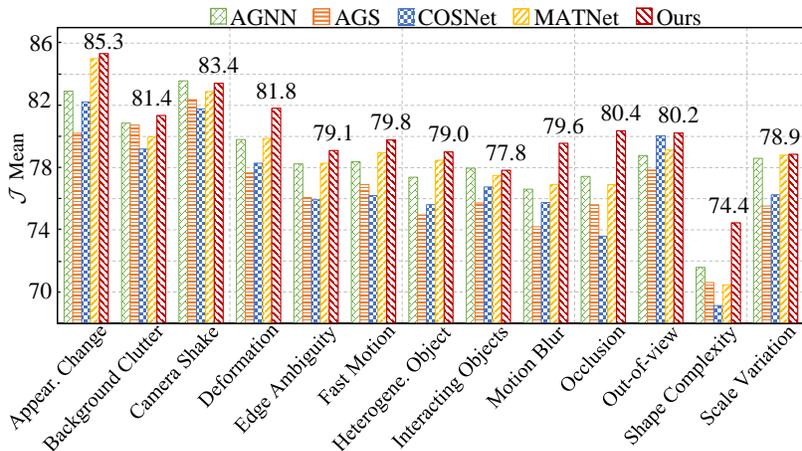
Figure 7: Performance on videos with various visual characteristics.

for them can model spatial pixel correlations better.

**Tradeoff parameter** $\alpha$. We vary the tradeoff parameter $\alpha$ in the loss function from 0.1 to 0.9 at an interval of 0.2, and the results are shown in Table 13. The results show that the performance of our method tends to rise up before 0.5 and saturates after the best value 0.5. This indicates that our method performs best when the model loss $\mathcal{L}_1$ and the knowledge distillation loss $\mathcal{L}_2$ contribute equally to the objective function.

Table 13: Tradeoff parameter $\alpha$ in the loss.

| $\alpha$ | $\overline{\mathcal{J}}$ | $\overline{\mathcal{F}}$ | $\overline{\mathcal{J}\&\mathcal{F}}$ |
|---|---|---|---|
| 0.1 | 78.3 | 80.2 | 79.3 |
| 0.3 | 80.1 | 81.5 | 80.8 |
| 0.5 | **82.4** | **84.3** | **83.4** |
| 0.7 | 81.9 | 82.8 | 82.4 |
| 0.9 | 81.8 | 83.0 | 82.4 |

**Various visual characteristics**. We show the LSTA performance on video data with 13 kinds of visual characteristics in Fig. 7. As depicted in this histogram, our model consistently performs better across a wide range of data characteristics. For example, LSTA is higher than the most competitive one by a large margin in several challenging scenarios, including background clutter, deformation, fast motion, motion blur, and occlusion. This provides solid evidence of the strong robustness of our model.

*4.6. Qualitative Results*

To give an intuitive view on the superiority of our model, we visualize segmentation results of randomly selected video frames from DAVIS2016 [26], YouTube-objects [28], and FBMS [24] in Fig. 8. As drawn in the figure, our model can give an accurate mask of the human-bicycle object regardless of the occluded tree (row 1), and successfully identify the local parts of objects, e.g., the feet of animals (cats in row 3 and dogs in row 5). This demonstrates that our model can well capture the moving object and the local pattern of object. In the meantime, Fig. 9 shows the masks of multiple objects in four video sequences randomly chosen from DAVIS2017, and the results validate that LSTA enjoys satisfying discriminative ability of distinct categories in the same scenario.

Moreover, we exhibit the visualized feature maps generated during the intermediate procedures of inference on a video from DAVIS2016 in Fig. 10. The first row shows the appearance feature map after Encoder, the second row shows the enhanced appearance feature map after Conv2D, the third and the fourth rows show the global and the local feature maps by passing LTM and STA, respectively. As vividly illustrated in these images, global feature maps are good at encoding the object contours, and local feature maps indeed play a similar role of optical flow in discovering the pattern of object motion.
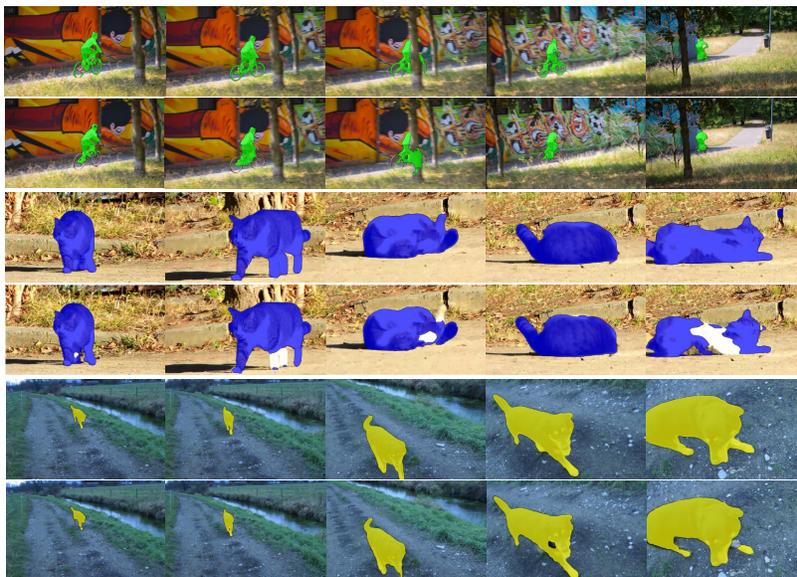
18

Figure 8: Segmentation results of LSTA (Row 1, 3, 5) and MATNet (row 2, 4, 6) on three randomly selected video from DAVIS2016 (row 1&2), YouTube-objects (row 3&4), and FBMS (row 5&6) data sets, respectively.
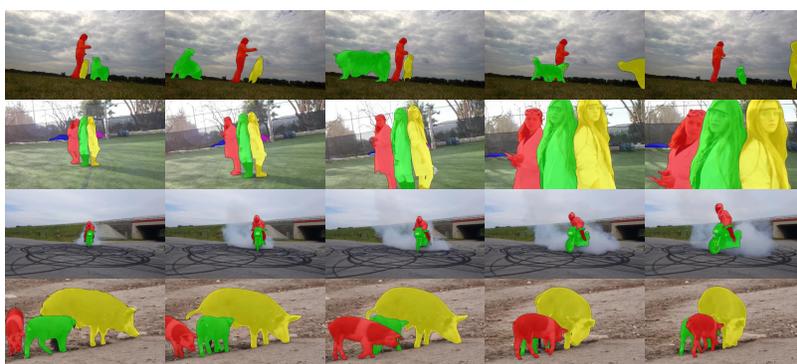


Figure 9: Segmentation results of multiple objects by LSTA on DAVIS2017.

In addition, we show some failure cases of our model on DAVIS2016 in Fig. 11. For primary objects, such as a dancing man in case (a) and a running car in case (b), there are noisy objects with similar appearance, e.g., the boy closest to the dancing man and another car in the right corner, which prevents our model from accurately segmenting the primary objects. This might be because our model is dependent of visual appearance, so one can resort to instance segmentation to alleviate the problem of similarity interference in appearance.

## 5. Conclusion

We have developed an end-to-end real-time unsupervised video object segmentation approach, named LSTA. It includes two primary blocks, i.e., Long Temporal Memory and Short Temporal Attention, which encode both long-range and short-range spatiotemporal pixel relations of the current frame and the past frames, respectively. The former LTM captures those constantly present objects from a global view, while the latter STA models the pattern of moving objects from a local view. Moreover, we have explored the performance of our method on several benchmark datasets. Both quantitative records and qualitative visualization results indicate the superiority of the proposed approach, including more promising segmentation masks, real-time inference speed, and robustness to some deformations or occlusions.
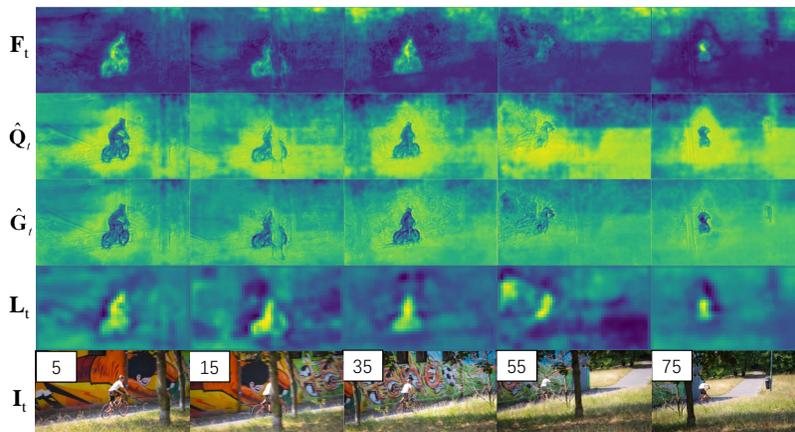
Figure 10: Feature map visualization of LSTA on DAVIS2016.
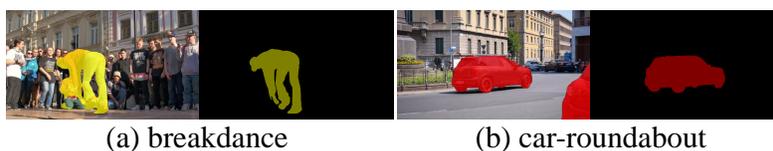


(a) breakdance      (b) car-roundabout

Figure 11: Failure cases on DAVIS2016. The former is the predicted mask by LSTA, the latter is the ground-truth mask.

There still exist some limitations in our LSTA method to be addressed in the future. 1) Due to the lacking of supervision frame, it is difficult to capture those small or tiny objects, and which can be solved by designing unsupervised VOS methods tailored for small objects. 2) It fails to handle the objects with occlusions, which often appear in real-world applications. Hence, a heuristic way is to adopt the video inpainting method to recover the occluded parts first. 3) Compared to the current vanilla method, it will be interesting to consider using additional knowledge, such as referring expressions and object detections, so as to further improve the performance in unsupervised setting.

## Acknowledgment

## References

[1] Ali Athar, Sabarinath Mahadevan, Aljosa Osep, Laura Leal-Taixé, and Bastian Leibe. STEm-Seg: Spatio-temporal embeddings for instance segmentation in videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 158–177, 2020.

[2] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4974–4983, 2019.

[3] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 801–818, 2018.

[4] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. In *The International Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

[5] Jingchun Cheng, Yi-Hsuan Tsai, Shengjin Wang, and Ming-Hsuan Yang. Segflow: Joint learning for video object segmentation and optical flow. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 686–695, 2017.

[6] Suhwan Cho, Minhyeok Lee, Seunghoon Lee, Chaewon Park, Donghyeong Kim, and Sangyoun Lee. Treating motion as option to reduce motion dependency in unsupervised video object segmentation. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 5129–5138, 2023.

[7] Mingqi Gao, Jungong Han, Feng Zheng, James J.Q. Yu, and Giovanni Montana. Video object segmentation using point-based memory network. *Pattern Recognition (PR)*, 134:109073, 2023.

[8] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 42(2):386–397, 2020.

[9] Ge-Peng Ji, Keren Fu, Zhe Wu, Deng-Ping Fan, Jianbing Shen, and Ling Shao. Full-duplex strategy for video object segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 4922–4933, 2021.

[10] Wei Ji, Jingjing Li, Qi Bi, Wenbo Li, and Li Cheng. Segment anything is not always perfect: an investigation of sam on different real-world applications. *arXiv*, abs/2304.05750, 2023.

[11] Choromanski Krzysztof, Likhosherstov Valerii, Dohan David, Song Xingyou, Gane Andreea, Sarlos Tamas, Hawkins Peter, Davis Jared, Mohiuddin Afroz, Kaiser Lukasz, Belanger David, Colwell Lucy, and Weller Adrian. Rethinking attention with performers. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.

[12] Meng Lan, Jing Zhang, and Zengmao Wang. Coherence-aware context aggregator for fast video object segmentation. *Pattern Recognition (PR)*, 136:109214, 2023.

[13] Youngjo Lee, Hongje Seong, and Euntai Kim. Iteratively selecting an easy reference frame makes unsupervised video object segmentation easier. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 1245–1253, 2022.

[14] Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3D semantic scene completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3351–3359, 2020.

[15] Mengze Li, Han Wang, Wenqiao Zhang, Jiaxu Miao, Zhou Zhao, Shengyu Zhang, Wei Ji, and Fei Wu. Winner: weakly-supervised hierarchical decomposition and alignment for spatio-temporal video grounding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 23090–23099, June 2023.

[16] Siyang Li, Bryan Seybold, Alexey Vorobyov, Xuejin Lei, and C.-C. Jay Kuo. Unsupervised video object segmentation with motion-based bilateral networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 215–231, 2018.

[17] Yu Li, Zhuoran Shen, and Ying Shan. Fast video object segmentation using the global context module. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 735–750, 2020.

[18] Chen Liang, Wenguan Wang, Tianfei Zhou, Jiaxu Miao, Yawei Luo, and Yi Yang. Local-global context aware transformer for language-guided video segmentation. *IEEE Transactions Pattern Analysis and Machine Intelligence (TPAMI)*, 45(8):10055–10069, 2023.

[19] Xiankai Lu, Wenguan Wang, Martin Danelljan, Tianfei Zhou, Jianbing Shen, and Luc Van Gool. Video object segmentation with episodic graph memory networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, volume 12348, pages 661–679, 2020.

[20] Xiankai Lu, Wenguan Wang, Chao Ma, Jianbing Shen, Ling Shao, and Fatih Porikli. See More, Know More: Unsupervised video object segmentation with co-attention siamese networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3623–3632, 2019.

[21] Jonathon Luiten, Idil Esen Zulfikar, and Bastian Leibe. UnOVOST: Unsupervised offline video object segmentation and tracking. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1989–1998, 2020.

[22] Sabarinath Mahadevan, Ali Athar, Aljosa Osep, Laura Leal-Taixé, Bastian Leibe, and Sebastian Hennen. Making a case for 3d convolutions for object segmentation in videos. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2020.

[23] Etienne Meunier, Anaïs Badoual, and Patrick Bouthemy. Em-driven unsupervised learning for efficient motion segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 45(4):4462–4473, 2023.

[24] Peter Ochs, Jitendra Malik, and Thomas Brox. Segmentation of moving objects by long term video analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 36(6):1187–1200, 2014.

[25] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 9226–9235, 2019.

[26] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 724–732, 2016.

[27] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017.

[28] Alessandro Prest, Christian Leistner, Javier Civera, Cordelia Schmid, and Vittorio Ferrari. Learning object class detectors from weakly annotated video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3282–3289, 2012.

[29] Sucheng Ren, Wenxi Liu, Yongtuo Liu, Haoxin Chen, Guoqiang Han, and Shengfeng He. Reciprocal transformations for unsupervised video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15455–15464, 2021.

[30] Sucheng Ren, Wenxi Liu, Yongtuo Liu, Haoxin Chen, Guoqiang Han, and Shengfeng He. Reciprocal transformations for unsupervised video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15455–15464, 2021.

[31] Andreas Robinson, Felix Jaremo Lawin, Martin Danelljan, Fahad Shahbaz Khan, Michael Felsberg, Felix Jaremo Lawin, Martin Danelljan, Fahad Shahbaz Khan, and Michael Felsberg. Learning fast and robust target models for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7406–7415, 2020.

[32] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 761–769, 2016.

[33] Hongmei Song, Wenguan Wang, Sanyuan Zhao, Jianbing Shen, and Kin-Man Lam. Pyramid dilated deeper con-vlstm for video salient object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 744–760, 2018.

[34] Jiadai Sun, Yuxin Mao, Yuchao Dai, Yiran Zhong, and Jianyuan Wang. Munet: Motion uncertainty-aware semi-supervised video object segmentation. *Pattern Recognition (PR)*, 138:109399, 2023.

[35] Mingjie Sun, Jimin Xiao, Eng Gee Lim, Yanchun Xie, and Jiashi Feng. Adaptive roi generation for video object segmentation using reinforcement learning. *Pattern Recognition (PR)*, 106:107465, 2020.

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *The International Conference on Neural Information Processing Systems (NeurIPS)*, pages 5999–6009, 2017.

[37] Carles Ventura, Miriam Bellver, Andreu Girbau, Amaia Salvador, Ferran Marques, and Xavier Giro-I-Nieto. RVOS: End-to-end recurrent network for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5272–5281, 2019.

[38] Paul Voigtlaender, Yuning Chai, Florian Schroff, Hartwig Adam, Bastian Leibe, and Liang Chieh Chen. Feelvos: Fast end-to-end embedding learning for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9473–9482, 2019.

[39] Haochen Wang, Xiaolong Jiang, Haibing Ren, Yao Hu, and Song Bai. Swiftnet: Real-time video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1296–1305, 2021.

[40] Wenguan Wang, Xiankai Lu, Jianbing Shen, David J Crandall, and Ling Shao. Zero-shot video object segmentation via attentive graph neural networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 9236–9245, 2019.

[41] Wenguan Wang, Jianbing Shen, Xiankai Lu, Steven C. H. Hoi, and Haibin Ling. Paying attention to video object pattern understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(7):2413–2428, 2021.

[42] Wenguan Wang, Jianbing Shen, Ruigang Yang, and Fatih Porikli. Saliency-aware video object segmentation. *IEEE Transactions Pattern Analysis and Machine Intelligence (TPAMI)*, 40(1):20–33, 2018.

[43] Lin Xi, Weihai Chen, Xingming Wu, Zhong Liu, and Zhengguo Li. Implicit motion-compensated network for unsupervised video object segmentation. *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, 32(9):6279–6292, 2022.

[44] Huaxin Xiao, Bingyi Kang, Yu Liu, Maojun Zhang, and Jiashi Feng. Online meta adaptation for fast video object segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 42(5):1205–1217, 2020.

[45] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtubevos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018.

[46] Shuangjie Xu, Daizong Liu, Linchao Bao, Wei Liu, and Pan Zhou. Mhp-vos: Multiple hypotheses propagation for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 314–323, 2019.

[47] Yanchao Yang, Brian Lai, and Stefano Soatto. Dystab: Unsupervised object segmentation via dynamic-static bootstrapping. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2826–2836, 2021.

[48] Zhao Yang, Qiang Wang, Luca Bertinetto, Weiming Hu, Song Bai, and Philip H S Torr. Anchor diffusion for unsupervised video object segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 931–940, 2019.

[49] Yingjie Yin, De Xu, Xingang Wang, and Lei Zhang. Agunet: Annotation-guided u-net for fast one-shot video object segmentation. *Pattern Recognition (PR)*, 110:107580, 2021.

[50] Yingjie Yin, De Xu, Xingang Wang, and Lei Zhang. Directional deep embedding and appearance learning for fast video object segmentation. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 33(8):3884–3894, 2022.

[51] Kaihua Zhang, Zicheng Zhao, Dong Liu, Qingshan Liu, and Bo Liu. Deep transport network for unsupervised video object segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 8761–8770, 2021.

[52] Lu Zhang, Jianming Zhang, Zhe Lin, Radomír Měch, Huchuan Lu, and You He. Unsupervised video object segmentation with joint hotspot tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 490–506, 2020.

[53] Zongji Zhao, Sanyuan Zhao, and Jianbing Shen. Real-time and light-weighted unsupervised video object segmentation network. *Pattern Recognition (PR)*, 120:108120, 2021.

[54] Tianfei Zhou, Jianwu Li, Xueyi Li, and Ling Shao. Target-aware object discovery and association for unsupervised video multi-object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6985–6994, 2021.

[55] Tianfei Zhou, Jianwu Li, Shunzhou Wang, Ran Tao, and Jianbing Shen. MATNet: Motion-attentive transition network for zero-shot video object segmentation. *IEEE Transactions on Image Processing (TIP)*, 29:8326–8338, 2020.

[56] Tianfei Zhou, Fatih Porikli, David J. Crandall, Luc Van Gool, and Wenguan Wang. A survey on deep learning technique for video segmentation. *IEEE Transactions Pattern Analysis and Machine Intelligence (TPAMI)*, 45(6):7099–7122, 2023.

[57] Yifeng Zhou, Xing Xu, Fumin Shen, Xiaofeng Zhu, and Heng Tao Shen. Flow-edge guided unsupervised video object segmentation. *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, 32(12):8116–8127, 2022.