

Recurrent Dynamic Embedding for Video Object Segmentation

Mingxing Li^{1*}, Li Hu^{2*}, Zhiwei Xiong^{1†}, Bang Zhang², Pan Pan², Dong Liu¹

¹University of Science and Technology of China

²Alibaba DAMO Academy, Alibaba Group

mxli@mail.ustc.edu.cn {zwxiong, dongeliu}@ustc.edu.cn

{hooks.hl, zhangbang.zb, panpan.pp}@alibaba-inc.com

Abstract

Space-time memory (STM) based video object segmentation (VOS) networks usually keep increasing memory bank every several frames, which shows excellent performance. However, 1) the hardware cannot withstand the ever-increasing memory requirements as the video length increases. 2) Storing lots of information inevitably introduces lots of noise, which is not conducive to reading the most important information from the memory bank. In this paper, we propose a Recurrent Dynamic Embedding (RDE) to build a memory bank of constant size. Specifically, we explicitly generate and update RDE by the proposed Spatio-temporal Aggregation Module (SAM), which exploits the cue of historical information. To avoid error accumulation owing to the recurrent usage of SAM, we propose an unbiased guidance loss during the training stage, which makes SAM more robust in long videos. Moreover, the predicted masks in the memory bank are inaccurate due to the inaccurate network inference, which affects the segmentation of the query frame. To address this problem, we design a novel self-correction strategy so that the network can repair the embeddings of masks with different qualities in the memory bank. Extensive experiments show our method achieves the best tradeoff between performance and speed. Code is available at <https://github.com/Limingxing00/RDE-VOS-CVPR2022>.

1. Introduction

Video object segmentation (VOS) is a fundamental task for video understanding, including lots of applications, such as autonomous driving and video editing. This work focuses on semi-supervised VOS setting. In this setting, given the instances annotation of the first frame, the VOS algorithms segment the instances in other frames.

Matching based networks [5, 13, 18, 20, 23–25, 30, 31,

* Equal contribution. † Corresponding author. This work was done during Mingxing Li’s internship at Alibaba.

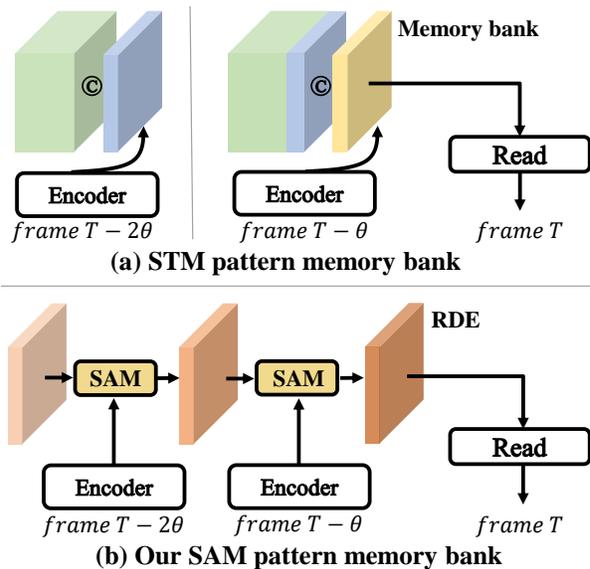


Figure 1. The inference pipelines of the segmentation of frame T . \textcircled{C} denotes concatenation. θ denotes the sampling interval for the update of the memory bank. (a) shows the network read the space-time memory (STM) pattern memory bank to segment frame T . As the length of videos increases, the STM pattern memory bank has an ever-increasing size. In (b), we update a recurrent dynamic embedding (RDE) to build a memory bank of the constant size, which is maintained by a spatio-temporal aggregation module (SAM).

39, 44] are popular for semi-supervised VOS. These networks have a memory bank mechanism, which encodes some frames into embeddings and stores those embeddings in the memory bank to assist the segmentation of the query frame. Some methods only use the embeddings of a limited number of frames, such as the ground-truth (GT) frame [14], the latest frame (for brevity, the latest frame of the query frame is abbreviated as the latest frame) [27] and both of them [20, 24, 39]. These methods do not make full use of historical frames in the video. STM based methods [5, 13, 18, 23, 25, 30, 31, 44] store the embeddings every

several (e.g., 5) frames in the STM pattern memory bank as shown in Figure 1(a). Although STM based methods utilize equal interval sampling to mine the historical information in the video, as the length of videos increases, the STM pattern memory bank has an ever-increasing size and inevitably introduces lots of noise. Exponential moving average (EMA) based methods [17, 19, 33] try to address the problems. The EMA based methods index some pixel embeddings from the embeddings of the query frame and the memory bank according to certain criteria and fuse these pixel embeddings in the EMA way. However, the EMA based methods have a strong limitation because of the direct summation operation (see details in Sec. 3.1).

In this paper, we address two problems. 1) How to build and update a memory bank of the constant size to effectively and efficiently store historical information? 2) Except for the GT frame, other masks are inaccurate owing to the inaccurate network inference, how to correct the poor embedding encoded from the inaccurate masks?

For problem 1, we propose a recurrent dynamic embedding (RDE) to provide a richer representation for VOS. As shown in Figure 1(b), to generate and update RDE, we propose a spatio-temporal aggregation module (SAM) to organize the cue of the historical information (previous RDE) and the embedding of the latest frame adaptively. SAM includes three parts: *extracting*, *enhancing* and *squeezing*. The *extracting* part is responsible for organizing the spatio-temporal relationship between previous RDE and the embedding of the latest frame. Then, the *enhancing* part reinforces the spatio-temporal relationship and the *squeezing* part aggregates and compresses the spatio-temporal information. We refer to the memory bank maintained by SAM as the SAM pattern memory bank.

One potential risk of the SAM pattern memory bank is the recurrent update of RDE may cause error accumulation. However, we have no GT for training the generated RDE directly. To tackle this problem, we propose to employ auxiliary supervision for the distribution of RDE. In the training process, we additionally build a STM pattern memory bank (see Figure 1(a)) to obtain the uncompressed information and its read results, which are used to estimate the distribution for RDE. Thus we design an unbiased guidance loss to control the approach degree of the two distributions. Relying on the unbiased guidance loss, the training of the network is more stable and has higher performance with no extra computation overhead for deployment.

For problem 2, we design a novel self-correction strategy, which enforces the network to repair the embeddings of masks with different qualities in the memory bank. Specifically, we first simulate different perturbed masks and then constrain the embeddings encoded by perturbed masks to be close to the embedding encoded by the GT mask with a mask consistency loss. The mask consistency loss enforce

the network to learn the self-correction ability for inaccurate masks in the embedding space during the training stage.

To investigate the effectiveness of the proposed methods, we conduct experiments on DAVIS 2017, DAVIS 2016 and YouTube-VOS 2019. The proposed method achieves state-of-the-art performance on DAVIS 2017 validation set (86.1% $\mathcal{J}\&\mathcal{F}$, 27 FPS), DAVIS 2017 test set (78.9% $\mathcal{J}\&\mathcal{F}$), DAVIS 2016 (91.6% $\mathcal{J}\&\mathcal{F}$, 35 FPS) and superior performance on YouTube-VOS 2019 (83.3% $\mathcal{J}\&\mathcal{F}$) without the multi-scale inference. Furthermore, we demonstrate the effectiveness of our method in the synthetic long video. For the synthetic long video, $\mathcal{J}\&\mathcal{F}$ and FPS of our method are almost unchanged as the length of the synthetic long video increases.

Our contributions can be summarized as follows:

- We propose an easy-to-extend recurrent dynamic embedding (RDE) to provide a richer representation for VOS compared with the embedding of the GT frame and the latest frame, which is maintained by the proposed spatio-temporal aggregation module (SAM).
- To avoid error accumulation owing to the recurrent usage of SAM, we propose an unbiased guidance loss during the training stage, which makes SAM more robust in long videos.
- Considering inaccurately predicted masks in the memory bank affect the segmentation performance due to the inaccurate network inference, we design a novel self-correction strategy, which enforces the network to learn the self-correction ability for inaccurate masks in the embedding space.
- Extensive experiments on several benchmarks and the synthetic long video show the effectiveness and superiority of our method.

2. Related Work

2.1. Semi-supervised VOS

Semi-supervised VOS mainly focuses on propagating the certain object mask of one frame. It can be roughly divided semi-supervised VOS into three categories: 1) On-line fine-tuning based methods [22, 36], which usually learn general segmentation features and fine-tune the network to the target video during the test time. 2) Propagation based methods [3, 21], which refine the target segmentation mask in a temporal label propagation way. 3) Matching based methods [6, 23, 25, 30, 44] which encode some frames into embeddings and store those embeddings in the memory bank to segment the query frame.

2.2. Matching based VOS network

STM [25] is a popular network in matching based methods, which constructs a continuously updated memory bank

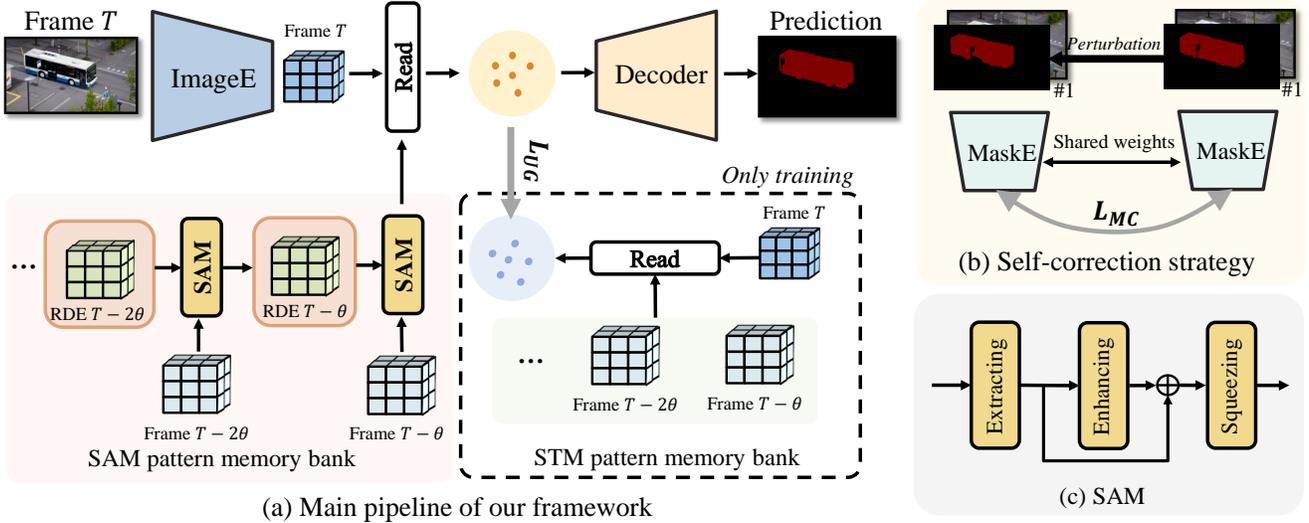


Figure 2. Illustrating the architectures: (a) The main pipeline of our framework. During the training stage, we maintain two individual memory banks which are updated in the STM pattern and our SAM pattern separately. During the inference, we only utilize our SAM pattern memory bank. θ denotes the sampling interval for the update of the memory bank. (b) Self-correction strategy. The proposed mask consistency loss L_{MC} enforces the mask encoder to learn the self-correction ability for the inaccurate masks. (c) The structure of SAM, which organizes the historical information and the embedding of the latest frame adaptively.

of historical frames. Compared with using limited frames (the GT frame [14] or the latest frame [27]), memory bank excavates the information of historical frames even more. Recently, matching based networks have received widespread attention. [5, 13, 20, 30] improve the readout operation of the memory bank, [37] applies a local attention with the help of the optical flow, [10] utilizes the global and instance embedding learning to address multi-objects VOS. Although these methods have achieved satisfactory performance, they ignore two key problems: 1) As the number of video frames increases, the hardware cannot afford the ever-increasing memory requirements. 2) Storing lots of information inevitably introduces lots of noise, which is not conducive to reading the most important information from the memory bank.

2.3. Efficient VOS network

The methods of efficient VOS usually belong to propagation based methods or matching based methods. SAT [3] is one of propagation based methods, which deals with each object as a tracklet and segments the object via two feedback loops. OSMN [39] is one of the matching based methods, which adopt the GT frame and the latest frame to guide the segmentation of the query frame with two modulators. Recently, the most popular inference setting for VOS is to save the feature embedding of historical frames every 5 frames (STM pattern). Some methods [17, 19, 33] try to use exponential moving average (EMA) to build a more efficient characterization to record the historical information. However, these methods only perform between the most

similar embeddings due to the direct summation operation (see details in Sec. 3.1), which is a strong limitation.

3. Method

3.1. Revisit Memory Bank Update with EMA

In STM [25], the image and mask are encoded into two embedding spaces, named *key* and *value*. In addition to the *key* and *value* of the GT frame and the latest frame, previous EMA based methods build an independent embedding, *IE*. Take *key* update as an example, let $\mathbf{k}_t^{IE}(p)$ denotes the *key* at time t and $\mathbf{k}^Q(q)$ denotes the *key* of the query frame Q , where p and q are the coordinate of the spatial position. [17, 19] utilize EMA to update the historical embedding $\mathbf{k}_{t-\theta}^{IE}(p)$ with the query embedding $\mathbf{k}^Q(q)$ by certain rules (see details in the supplementary material). The new embedding $\mathbf{k}_t^{IE}(p)$ in the memory bank can be formulated as follows:

$$\mathbf{k}_t^{IE}(p) = (1 - \lambda)\mathbf{k}^Q(q) + \lambda\mathbf{k}_{t-\theta}^{IE}(p) \quad (1)$$

where λ is a hyper-parameter to control the update strength and θ denotes the update interval. We argue that EMA based methods have a strong limitation, in which the two additional items in Eq. 1 must be similar in the parameter space because of the summation operation. Thus these methods [17, 19] index the most similar embeddings to update. Our method associates the embeddings to update the extra embedding adaptively.

3.2. Framework Overview

Encoders. The main pipeline of our framework is illustrated in Fig 2(a). For a query frame of size $H \times W$, the image encoder *ImageE* is responsible for extracting image features. We also adopt a mask encoder *MaskE* to encode a certain frame and its mask to store into the memory bank. Both the encoders adopt ResNet-50 [12] as the backbone and use two simple projection heads following STM [25] to obtain two embeddings, *key* $\mathbf{k} \in \mathbb{R}^{C_k \times \frac{H}{16} \times \frac{W}{16}}$ and *value* $\mathbf{v} \in \mathbb{R}^{C_v \times \frac{H}{16} \times \frac{W}{16}}$. Here C_k and C_v are the numbers of the channel dimension ($C_k = 64$, $C_v = 512$ in our experiments).

Memory Reading and Decoder. Following STCN [6], for the SAM pattern memory bank m at time t , we keep target-agnostic key \mathbf{k}_t^m and target-specific value $\mathbf{v}_{t,i}^m$, where i denotes the i -th object. For the similarity $\mathbf{S}(p, q)$ between the key from the SAM pattern memory bank $\mathbf{k}_t^m(p)$ and the key of the query frame $\mathbf{k}_t^Q(q)$, we perform negative squared Euclidean distance, which can be formulated as

$$\mathbf{S}(p, q) = -\|\mathbf{k}_t^m(p) - \mathbf{k}_t^Q(q)\|_2^2 \quad (2)$$

where p and q are the coordinate of the spatial position of $\mathbf{k}_t^m(p)$ and $\mathbf{k}_t^Q(q)$ separately. And the softmax operation is applied on the spatial dimension for similarity \mathbf{S} to obtain the softmax-normalized affinity matrix \mathbf{W} , $\mathbf{W} = \text{softmax}(\mathbf{S})$. Relying on \mathbf{W} , the readout feature $\mathbf{v}_{t,i}^{m \rightarrow Q}$ of the i -th object from the SAM memory bank can be obtained by the matrix multiplication \odot :

$$\mathbf{v}_{t,i}^{m \rightarrow Q} = \mathbf{W} \odot \mathbf{v}_{t,i}^m. \quad (3)$$

The readout feature $\mathbf{v}_{t,i}^{m \rightarrow Q}$ concatenates with the value of the query frame to pass through the light-weight decoder described in [6] to get the segmentation results $\tilde{\mathbf{y}}_{t,i}^m$ of the i -th object at frame t . Similar to the SAM pattern memory bank, we concatenate the readout feature $\mathbf{v}_{t,i}^{M \rightarrow Q}$ from the STM pattern memory bank M with the value of the query frame to obtain the segmentation results $\tilde{\mathbf{y}}_{t,i}^M$ of the i -th object at frame t .

3.3. SAM Pattern Memory Bank

The main challenge of keeping the size of the memory bank constant is how to select the most useful information. The STM pattern memory bank can store the historical information losslessly, but has ever-increasing size and inevitably introduces lots of noise. In our design, we build a SAM pattern memory bank to address the challenge. During the training stage, the STM and SAM pattern memory banks are maintained at the same time. During the inference, we only use the SAM pattern memory

bank, which can keep the size of the memory bank constant. Specifically, the STM pattern memory bank M includes $\{\mathbf{k}_t^M, \mathbf{v}_{t,i}^M\}$, while the SAM pattern memory bank m includes $\{\mathbf{k}_t^m, \mathbf{v}_{t,i}^m\}$.

Recurrent Dynamic Embedding. We find the embedding of the latest frame changes over time, providing more useful information for the segmentation of the query frame but lacking the use of historical information. We propose a recurrent dynamic embedding (RDE) in the memory bank to fuse the cue of the historical information with the embedding of the latest frame to provide a richer representation for VOS. We denote the RDE embedding at time t as $\{\mathbf{k}_t^{RDE}, \mathbf{v}_{t,i}^{RDE}\} \in \{\mathbf{k}_t^m, \mathbf{v}_{t,i}^m\}$.

Spatio-temporal Aggregation Module. To generate and update RDE, we propose a spatio-temporal aggregation module (SAM), which exploits the cue of historical information. SAM includes three parts: *extracting*, *enhancing* and *squeezing* as shown in Figure 2(c). The *extracting* part is responsible for organizing the spatio-temporal relationship between the embedding of previous RDE $\{\mathbf{k}_{t-\theta}^{RDE}, \mathbf{v}_{t-\theta,i}^{RDE}\}$ (θ denotes the sampling interval) and the embeddings of the latest frame $\{\mathbf{k}_t^L, \mathbf{v}_{t,i}^L\}$. First, we concatenate previous RDE $\{\mathbf{k}_{t-\theta}^{RDE}, \mathbf{v}_{t-\theta,i}^{RDE}\}$ and the embedding of the latest frame $\{\mathbf{k}_t^L, \mathbf{v}_{t,i}^L\}$ to obtain the feature \mathbf{x} . Take \mathbf{k}_t^{RDE} update as an example,

$$\mathbf{x} = \text{Cat}(\mathbf{k}_{t-\theta}^{RDE}, \mathbf{k}_t^L), \mathbf{x} \in \mathbb{R}^{C_k \times 2 \times \frac{H}{16} \times \frac{W}{16}} \quad (4)$$

where *Cat* denotes concatenation operation in the time dimension. Inspired by self-attention mechanism [35], in the *extracting* part, we organize the spatio-temporal relationship between previous RDE $\mathbf{k}_{t-\theta}^{RDE}$ and the embedding of the latest frame \mathbf{k}_t^L to obtain the aggregation feature \mathbf{x}_{agg} ,

$$\mathbf{x}_{agg} = \frac{1}{C(\mathbf{x})} \omega(\mathbf{x})^T \phi(\mathbf{x} \downarrow) g(\mathbf{x} \downarrow). \quad (5)$$

$C(\mathbf{x})$ is a normalization factor, which presents the total number of the spatial position of \mathbf{x} . The function ω , ϕ and g are $1 \times 1 \times 1$ convolution in our implementation. $\mathbf{x} \downarrow$ denotes \mathbf{x} processed by the max-pooling operation (no down sampling on the time axis), which can decrease the computational complexity.

Relying on the aggregation feature \mathbf{x}_{agg} , in the *enhancing* part, we enhance \mathbf{x}_{agg} in the form of residuals by atrous spatial pyramid pooling (ASPP) [2]. Finally, in the *squeezing* part, we compress the enhanced feature by a simple $2 \times 3 \times 3$ convolution, which is denoted as *Squeeze* function. The formula can be expressed as

$$\mathbf{k}_t^{RDE} = \text{Squeeze}(\mathbf{x}_{agg} + \text{ASPP}(\mathbf{x}_{agg})). \quad (6)$$

previous RDE and the embedding of the latest frame adaptively fuse and maintain the constant size for the memory bank, where the key mapping is $\mathbb{R}^{C_k \times 2 \times \frac{H}{16} \times \frac{W}{16}} \rightarrow \mathbb{R}^{C_k \times 1 \times \frac{H}{16} \times \frac{W}{16}}$. For the multiple objects, we concatenate the object dimension to the batch dimension like the implementation of STM [25]. For the key and value of RDE, we maintain two different SAMs respectively.

Unbiased Guidance Loss. One potential risk of the SAM pattern memory bank is the update of RDE may cause error accumulation, especially when it is used repeatedly. Another problem is that the key and value of RDE are generated separately by two different SAMs, the distribution of them is difficult to directly define. Suppose the update process of the STM pattern memory bank is a good teacher, the estimated distribution read from the SAM pattern memory bank ought to approach the estimated distribution read from the STM pattern memory bank. Thus, during the training stage, we maintain two individual memory banks for the segmentation of the query frame, which is updated in the STM and SAM patterns separately. We propose an unbiased guidance loss L_{UG} , which controls the distribution of the readout feature from the SAM pattern memory bank $\mathbf{v}_{t,i}^{m \rightarrow Q}$ to approach the distribution of the readout feature from the STM pattern memory bank $\mathbf{v}_{t,i}^{M \rightarrow Q}$. The unbiased guidance loss L_{UG} is computed as follows:

$$L_{UG} = \sum_i KL(\mathbf{v}_{t,i}^{M \rightarrow Q} || \mathbf{v}_{t,i}^{m \rightarrow Q}). \quad (7)$$

KL function denotes Kullback–Leibler (KL) divergence, which is a non-symmetric measure of the difference between two distributions.

Self-correction Strategy. Considering the quality of the mask in the memory bank affects the segmentation of the query frame, we propose a mask consistency loss L_{MC} to constrain the consistency of the embedding of masks of different qualities and the embedding of the GT mask. We first obtain the key \mathbf{k}_1 and value $\mathbf{v}_{1,i}$ of the first frame. And we perform perturbation transform such as the random dilation and eroding on the first frame to obtain the perturbed key $\check{\mathbf{k}}_1$ and perturbed value $\check{\mathbf{v}}_{1,i}$. The mask consistency loss L_{MC} can be calculated by

$$L_{MC} = KL(\mathbf{k}_1 || \check{\mathbf{k}}_1) + \sum_i KL(\mathbf{v}_{1,i} || \check{\mathbf{v}}_{1,i}) \quad (8)$$

where KL function denotes KL divergence.

Overall Loss Functions. During the training stage, we sample 5 frames. Inspired by the slowfast network [9], we utilize the SAM pattern memory bank to segment the third and fifth frames to handle different rate of videos. Besides,

we utilize STM pattern memory bank to segment the second and fourth frames for the training stability. We use bootstrapped cross entropy (BCE) following [5] to supervise the final segmentation results, which is computed as follows:

$$L_{Seg} = \frac{1}{2} \left(\sum_i \sum_{t=2,4} \underbrace{BCE(\tilde{\mathbf{y}}_{t,i}^M, \mathbf{y}_{t,i})}_{STM \text{ pattern item}} + \sum_i \sum_{t=3,5} \underbrace{BCE(\tilde{\mathbf{y}}_{t,i}^m, \mathbf{y}_{t,i})}_{SAM \text{ pattern item}} \right) \quad (9)$$

where $\tilde{\mathbf{y}}_{t,i}^M$ and $\tilde{\mathbf{y}}_{t,i}^m$ denote the segmentation results read from the STM pattern memory bank and the SAM pattern memory bank separately. $\mathbf{y}_{t,i}$ denotes the GT mask of the i -th object at frame t . The overall loss function is computed as follows:

$$Loss = L_{Seg} + \mathbb{1}[t = 3, 5] \mu L_{UG} + \gamma L_{MC} \quad (10)$$

where μ and γ are hyper-parameters to control the strength. We set $\mu = 10$ and $\gamma = 10$ in our experiments. $\mathbb{1}[\cdot]$ is the indicator function.

Inference Strategy. As shown in Figure 2(a), during the inference, we employ SAM recurrently to update RDE. Specifically, in a video of any length, SAM inputs previous RDE at time $t - \theta$ and the embeddings of the latest frame at time t to generate RDE at time t , where θ is the sampling interval. The new RDE is stored in the SAM pattern memory bank to assist the segmentation of the query frame and the old RDE is discarded.

4. Experiments

4.1. Datasets and Metrics

DAVIS. DAVIS 2016 [28] is a popular benchmark for video single object segmentation, whose validation set includes 20 videos. DAVIS 2017 [29] is a popular benchmark for video multiple objects segmentation, whose validation set and test set are 30 densely annotated videos.

YouTube-VOS. YouTube-VOS 2019 [38] is a large-scale benchmark for multi-object video segmentation, providing 3,471 videos for the training (65 categories) and 507 videos for the validation. There are additional 26 unseen categories in the validation set for evaluating the generalization.

Metrics. For the DAVIS datasets, we use the region similarity \mathcal{J} , the contour accuracy \mathcal{F} and their average $\mathcal{J\&F}$ to evaluate the segmentation results. For YouTube-VOS 2019, we follow the official evaluation server to report \mathcal{J} and \mathcal{F} of the seen and unseen categories, and the average of them.

4.2. Implementation Details

Training Stages. Following STCN [6], we first train the network equipped with the STM pattern memory bank on the static datasets [4, 16, 32, 34, 43] with 75k iterations and batch size of 64. The static images are processed by synthetic deformations like STM [25]. Secondly, we train the network equipped with the SAM and STM pattern memory banks on BL30K [1, 7] proposed in [5] with 500k iterations and batch size of 8. Finally, we fine-tune the network equipped with the SAM and STM pattern memory banks on YouTube-VOS and DAVIS 2017 with 75k iterations and batch size of 16 (main stage). BatchNorm layers are frozen during the training stage following [25].

Training Details. We adopt four 16 GB Tesla V100 GPUs to implement Pytorch. All networks are optimized by Adam optimizer [15]. We pretrain the network on the static datasets and BL30K with the initial learning rate of $2e-5$ and $1e-5$. And we fine-tune the network on the main stage with the initial learning rate of $2e-5$. The data augmentation is the same as STCN [6]. Besides, we sample 3 frames in the first pre-training stage and 5 frames in other stages.

Inference Details. During inference, we only use the SAM pattern memory bank. Specifically, in addition to maintaining our RDE by SAM, we sample the embedding of the latest frame and two repeated embedding of the GT frame. This setting is to keep a sampling balance between the accurate template information (GT frame) and dynamic information (latest frame or our RDE). We use top-k filters [5] $k = 40$ on all datasets. The sampling interval θ is set to 3 on all DAVIS datasets and 4 on YouTube-VOS 2019.

4.3. Compare with the State-of-the-art Methods

We denote the memory bank of the constant size during the inference as **Constant Cost (CC)**. As the video length increases during the inference, the CC methods can maintain a relatively stable speed and constant requirements of the memory. For brevity, our Recurrent Dynamic Embedding for VOS method is denoted as **RDE-VOS**.

DAVIS. We compare the proposed method with previous state-of-the-art methods for VOS on the DAVIS 2017 validation set, DAVIS 2017 test set and DAVIS 2016 validation set. On the DAVIS 2017 validation set, as shown in Table 1, our method even outperforms STCN [6] by 0.7% for $\mathcal{J}\&\mathcal{F}$ and runs about 35% faster (27 vs 20.2 FPS). Compared with SwiftNet [33], our method suppresses it by 5% for $\mathcal{J}\&\mathcal{F}$ and has a slight advantage for the speed (+2 FPS). On the DAVIS 2017 test set, as shown in Table 2, our method still has great advantages. On DAVIS 2016 validation set, as shown in Table 3, our method outperforms CC

Method	CC	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	FPS
STM [†] [25]	×	81.8	79.2	84.3	10.2
KMN [†] [30]	×	82.8	80.0	85.6	<8.4
JOINT [†] [23]	×	83.5	80.8	86.2	4.0
LCM [†] [13]	×	83.5	80.5	86.5	<8.5
RMNet [†] [37]	×	83.5	81.0	86.0	<11.9
MiVOS ^{†*} [5]	×	84.5	81.7	87.4	11.2
HMMN [†] [31]	×	84.7	81.9	87.5	<10.0
STCN ^{†*} [6]	×	85.3	82.0	88.6	20.2
GCNet [17]	✓	71.4	69.3	73.5	<25.0
Liang <i>et al.</i> [19]	✓	74.6	73.0	76.1	4.0
G-FRTM [†] [26]	✓	76.4	-	-	18.2
PReMVOS [21]	✓	77.8	73.9	81.7	0.01
SwiftNet [†] [33]	✓	81.1	78.3	83.9	<25.0
SST [†] [8]	✓	82.5	79.9	85.1	-
Ge <i>et al.</i> [†] [10]	✓	82.7	80.2	85.3	6.7
RDE-VOS[†]	✓	84.2	80.8	87.5	27.0
RDE-VOS^{†*}	✓	86.1	82.1	90.0	27.0

Table 1. Results on the DAVIS 2017 validation set. CC denotes constant cost during the inference. [†] indicates YouTube-VOS [38] is added during the training stage. * denotes BL30K [5] is added during the training stage.

Method	CC	600p	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
STM [†] [25]	×	✓	72.2	69.3	75.2
KMN [†] [30]	×	✓	77.2	74.1	80.3
RMNet [†] [37]	×	×	75.0	71.9	78.1
Ge <i>et al.</i> [†] [10]	×	×	75.2	72.0	78.3
STCN ^{†*} [6]	×	×	77.8	74.3	81.3
MiVOS ^{†*} [5]	×	×	78.6	74.9	82.2
CFBI [†] [40]	✓	×	74.8	71.1	78.5
Ge <i>et al.</i> [†] [10]	✓	×	75.2	72.0	78.3
CFBI+ [†] [41]	✓	×	75.6	71.6	79.6
RDE-VOS[†]	✓	×	77.4	73.6	81.2
RDE-VOS^{†*}	✓	×	78.9	74.9	82.9

Table 2. Results on the DAVIS 2017 test set. 600p denotes evaluating on 600p resolution.

method SwiftNet [33] by 1.2% for $\mathcal{J}\&\mathcal{F}$ and runs about 40% faster (35 vs 25 FPS). Compared with STCN [6], our method is 30% faster while $\mathcal{J}\&\mathcal{F}$ is almost unchanged (-0.1%). We also show the qualitative result of the validation set of DAVIS 2017 in Figure 4. More qualitative results can be found in the supplementary material.

YouTube-VOS. On a large-scale YouTube-VOS 2019 validation set, we compare our method with recent state-of-the-art methods in Table 4. Although our method does not surpass STCN on YouTube-VOS 2019 validation set, it still surpasses other state-of-the-art methods, regardless of whether BL30K is added.

Method	CC	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}	FPS
RMNet [†] [37]	×	88.8	88.9	88.7	11.9
STM [†] [25]	×	89.3	88.7	89.9	6.3
KMN [†] [30]	×	90.5	89.5	91.5	8.4
LCM [†] [13]	×	90.7	89.9	91.4	8.5
HMMN [†] [31]	×	90.8	89.6	92.0	10.0
MiVOS ^{†*} [5]	×	91.0	89.7	92.4	16.9
STCN ^{†*} [6]	×	91.7	90.4	93.0	26.9
GCNet [17]	✓	86.6	87.6	85.7	25.0
CFBI+ [†] [41]	✓	89.9	88.7	91.1	5.9
SwiftNet [†] [33]	✓	90.4	90.5	90.3	25.0
RDE-VOS[†]	✓	91.1	89.7	92.5	35.0
RDE-VOS^{†*}	✓	91.6	90.0	93.2	35.0

Table 3. Results on the DAVIS 2016 validation set. CC denotes constant cost during the inference.

Method	CC	Overall	\mathcal{J}_{seen}	\mathcal{F}_{seen}	\mathcal{J}_{unseen}	\mathcal{F}_{unseen}
STM [†] [25]	×	79.2	79.6	83.6	73.0	80.6
MiVOS ^{†*} [5]	×	82.4	80.6	84.7	78.2	85.9
STCN ^{†*} [6]	×	84.2	82.6	87.0	79.4	87.7
CFBI [†] [40]	✓	81.0	80.6	85.1	75.2	83.0
SST [†] [8]	✓	81.8	80.9	-	76.6	-
RDE-VOS[†]	✓	81.9	81.1	85.5	76.2	84.8
RDE-VOS^{†*}	✓	83.3	81.9	86.3	78.0	86.9

Table 4. Results on the YouTube-VOS 2019 validation set.

Synthetic Long Video. Recently, the popular benchmarks include short video clips. For example, DAVIS 2017 only has 67 frames per video clip on average. However, many practical applications need to handle more frames. Compared with STCN [6], we demonstrate the effectiveness of our method in the scene where includes more frames. Take a DAVIS 2017 case “cows” (the basic length is 104) as an example, exerting video forward and backward as a basic unit, we repeatedly sample multiple basic units to synthesize a long video. This synthesis method ensures each frame contains GT and the adjacent frames have smooth transitions. As shown in Figure 3, as the length of the synthetic long video increases, the performance and speed of our method is almost unaffected, while the performance and speed of STCN obviously decrease. Here we do not change any hyper-parameter compared with the setting on the DAVIS datasets. Besides, we utilize the official code of STCN and minimize the sampling interval to 60 frames under maximizing the usage of the GPU memory. All input data is stored in the CPU and inferred on one GPU.

Inference Time. We evaluate the inference time on one Tesla V100 GPU with full floating point precision. On the validation set of DAVIS 2017 and DAVIS 2016, as shown in Table 1 and 3, our method has a great advantage in speed compared with STCN (27 vs 20.2 FPS on DAVIS 2017 and 35 vs 26.9 FPS on DAVIS 2016).

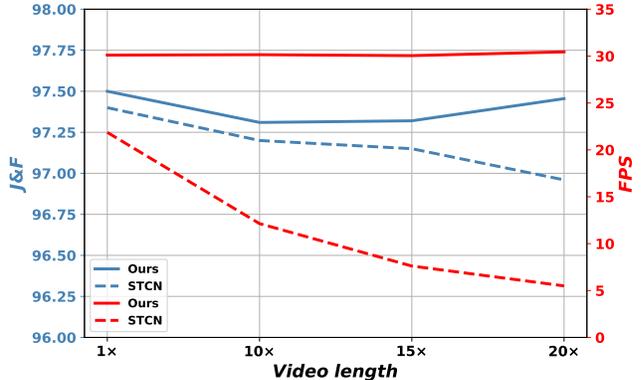


Figure 3. $\mathcal{J}\&\mathcal{F}$ and FPS of our method and STCN [6] on the synthetic long video. Note different colored lines refer to different metrics. When the length of the synthetic long video is 1, 10, 15 and 20 times of the original, $\mathcal{J}\&\mathcal{F}$ and FPS of our method are almost unchanged. However, both $\mathcal{J}\&\mathcal{F}$ and FPS of STCN have an obvious reduction.

4.4. Ablation Study

Dataset Setting. We compare the results whether to adopt BL30K [5] in Table 1, 2, 3 and 4. Without the BL30K pre-training, our method has superior performance on all datasets with a higher speed compared with other state-of-the-art methods. After adding the BL30K pre-training, our method has a stable improvement on all datasets.

Inference Setting. Table 5 shows different inference strategies adopting the memory bank. Compared with only using the embedding from the first frame or the latest frame, only using our RDE has the best performance of 81.8% for $\mathcal{J}\&\mathcal{F}$. Besides, based on using the embeddings from the first frame, the latest frame, and both of them, adding our RDE can further improve $\mathcal{J}\&\mathcal{F}$ by 13.7%, 1.8% and 0.8% separately. Based on using RDE and the embedding of both the first frame and the latest frame, we explore the sampling balance of the accurate template information (GT frame) and dynamic information (latest frame or our RDE). We find additionally sample the embedding of the GT frame to keep the sampling balance between the two types of information can further improve $\mathcal{J}\&\mathcal{F}$ by 0.7%. We use this strategy in all experiments unless otherwise specified. We also show the ablation of different sampling intervals θ , where the sampling interval of 3 provides the best result.

Loss Function Setting. In Table 6, we perform ablation of different loss functions without the BL30K [5] pre-training. Both our proposed L_{MC} and L_{UG} can improve the performance to different degrees, and their combination can maximize the performance (+1.7% $\mathcal{J}\&\mathcal{F}$). Besides, although we do not use the STM pattern memory bank during the inference, we find supervising the segmentation results

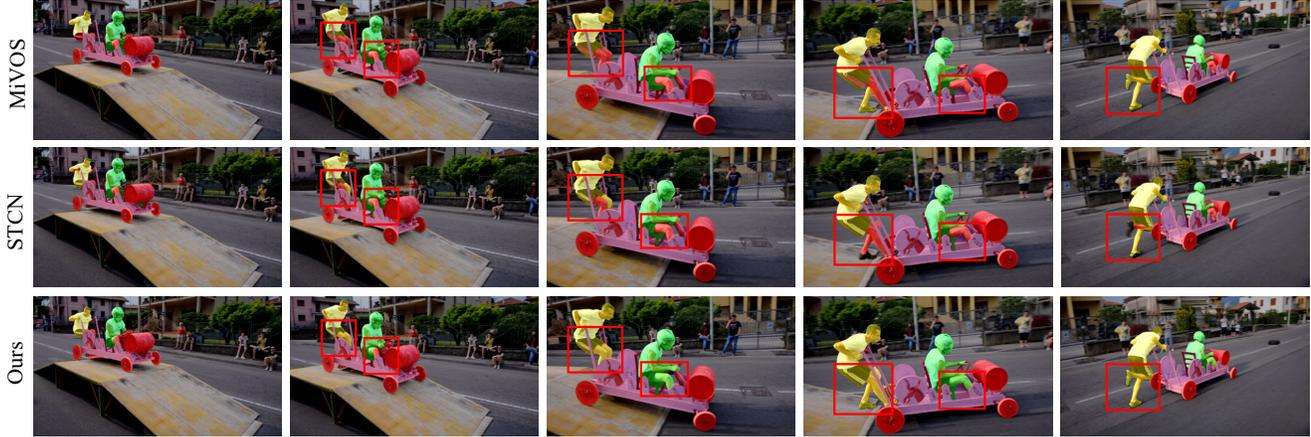


Figure 4. Qualitative results on the DAVIS 2017 validation set. We compare MiVOS [5] and STCN [6] under the challenging scale and deformation case, and our method has a notable improvement.

Variants	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
Strategy permutation			
RDE	81.8	78.0	85.7
First frame	71.6	67.8	75.4
First frame & RDE	85.3	81.6	89.0
Latest frame	80.4	76.9	83.8
Latest frame & RDE	82.2	78.4	86.0
First frame & latest frame	84.6	81.0	88.2
F & L & RDE	85.4	81.6	89.2
First frame $\times 2$ & latest frame	85.1	81.5	88.7
First frame & latest frame $\times 2$	84.0	80.4	87.6
2F & L & RDE	86.1	82.1	90.0
Sampling interval θ			
2F & L & RDE ($\theta = 2$)	85.1	81.4	88.9
2F & L & RDE ($\theta = 3$)	86.1	82.1	90.0
2F & L & RDE ($\theta = 4$)	85.1	81.5	88.8
2F & L & RDE ($\theta = 5$)	84.2	80.5	87.9

Table 5. Ablation of inference strategies on DAVIS 2017 validation set. F & L & RDE represents first frame, latest frame and RDE. 2F represents we sample the embeddings of the GT frame twice in order to keep balance of the accurate template information and dynamic information, which is used in all experiments unless otherwise specified.

	Ablation Settings	$\mathcal{J}\&\mathcal{F}$	\mathcal{J}	\mathcal{F}
Loss	w/o L_{MC}	83.7	80.5	86.9
	w/o L_{UG}	82.9	79.5	86.4
	w/o L_{MC} & L_{UG}	82.5	79.1	86.0
	L_{Seg} w/o STM pattern item	83.0	79.4	86.6
	Full	84.2	80.8	87.5

Table 6. Ablation of different loss functions without the BL30K [5] pre-training.

of the STM pattern item in Eq. 9 can assist the training of the SAM pattern (+1.2% $\mathcal{J}\&\mathcal{F}$).

4.5. Limitations.

During the inference, we set the sampling interval of the RDE update to θ . This simple setting is easy to plug in other matching based VOS methods. We fix the sampling interval $\theta = 3$ on the DAVIS datasets and achieve the new state-of-the-art performance. We increase the sampling interval on YouTube-VOS by 1 to fit the motion pattern on YouTube-VOS. A better solution for future work is to use a learnable discriminator [11] or gate mechanism [42] to adaptively control the update interval of SAM, which can handle different scenes better.

5. Conclusion

In this paper, we explore how to build and update a memory bank of the constant size to maximize the segmentation performance of the query frame. The key insight is we propose a recurrent dynamic embedding (RDE) to provide a richer representation for VOS compared with the embeddings of the GT frame and the latest frame. To generate and update RDE, we propose a novel spatio-temporal aggregation module (SAM), which organizes the cue of the historical information and the embedding of the latest frame adaptively. To avoid error accumulation owing to the recurrent usage of SAM, we propose an unbiased guidance loss during the training stage, which makes SAM more robust in long videos. Besides, we design a novel self-correction strategy so that the network can encode and self-repair the embeddings of masks with different qualities.

Acknowledgement. We acknowledge funding from National Key R&D Program of China under Grant 2017YFA0700800, National Natural Science Foundation of China under Grants 61931014, 62131003 and 62021001, and Fundamental Research Funds for the Central Universities under No. WK3490000006.

References

- [1] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. [6](#)
- [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. [4](#)
- [3] Xi Chen, Zuoxin Li, Ye Yuan, Gang Yu, Jianxin Shen, and Donglian Qi. State-aware tracker for real-time video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9384–9393, 2020. [2, 3](#)
- [4] Ho Kei Cheng, Jihoon Chung, Yu-Wing Tai, and Chi-Keung Tang. Cascadepsp: toward class-agnostic and very high-resolution segmentation via global and local refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8890–8899, 2020. [6](#)
- [5] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5559–5568, 2021. [1, 3, 5, 6, 7, 8](#)
- [6] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. *arXiv preprint arXiv:2106.05210*, 2021. [2, 4, 6, 7, 8](#)
- [7] Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad Elbadrawy, Ahsan Lodhi, and Harinandan Katam. Blenderproc. *arXiv preprint arXiv:1911.01911*, 2019. [6](#)
- [8] Brendan Duke, Abdalla Ahmed, Christian Wolf, Parham Aarabi, and Graham W Taylor. Sstvos: Sparse spatiotemporal transformers for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5912–5921, 2021. [6, 7](#)
- [9] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. [5](#)
- [10] Wenbin Ge, Xiankai Lu, and Jianbing Shen. Video object segmentation using global and instance embedding learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16836–16845, 2021. [3, 6](#)
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [8](#)
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [4](#)
- [13] Li Hu, Peng Zhang, Bang Zhang, Pan Pan, Yinghui Xu, and Rong Jin. Learning position and target consistency for memory-based video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4144–4154, 2021. [1, 3, 6, 7](#)
- [14] Yuan-Ting Hu, Jia-Bin Huang, and Alexander G Schwing. Videomatch: Matching based video object segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 54–70, 2018. [1, 3](#)
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [6](#)
- [16] Xiang Li, Tianhan Wei, Yau Pun Chen, Yu-Wing Tai, and Chi-Keung Tang. Fss-1000: A 1000-class dataset for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2869–2878, 2020. [6](#)
- [17] Yu Li, Zhuoran Shen, and Ying Shan. Fast video object segmentation using the global context module. In *European Conference on Computer Vision*, pages 735–750. Springer, 2020. [2, 3, 6, 7](#)
- [18] Shuxian Liang, Xu Shen, Jianqiang Huang, and Xian-Sheng Hua. Video object segmentation with dynamic memory networks and adaptive object alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8065–8074, 2021. [1](#)
- [19] Yongqing Liang, Xin Li, Navid Jafari, and Qin Chen. Video object segmentation with adaptive feature bank and uncertain-region refinement. *arXiv preprint arXiv:2010.07958*, 2020. [2, 3, 6](#)
- [20] Xiankai Lu, Wenguan Wang, Martin Danelljan, Tianfei Zhou, Jianbing Shen, and Luc Van Gool. Video object segmentation with episodic graph memory networks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 661–679. Springer, 2020. [1, 3](#)
- [21] Jonathon Luiten, Paul Voigtlaender, and Bastian Leibe. Pre-mvos: Proposal-generation, refinement and merging for video object segmentation. In *Asian Conference on Computer Vision*, pages 565–580. Springer, 2018. [2, 6](#)
- [22] K-K Maninis, Sergi Caelles, Yuhua Chen, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. Video object segmentation without temporal information. *IEEE transactions on pattern analysis and machine intelligence*, 41(6):1515–1530, 2018. [2](#)
- [23] Yunyao Mao, Ning Wang, Wengang Zhou, and Houqiang Li. Joint inductive and transductive learning for video object segmentation. *arXiv preprint arXiv:2108.03679*, 2021. [1, 2, 6](#)
- [24] Seoung Wug Oh, Joon-Young Lee, Kalyan Sunkavalli, and Seon Joo Kim. Fast video object segmentation by reference-guided mask propagation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7376–7385, 2018. [1](#)
- [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/CVF International*

- Conference on Computer Vision*, pages 9226–9235, 2019. 1, 2, 3, 4, 5, 6, 7
- [26] Hyojin Park, Jayeon Yoo, Seohyeong Jeong, Ganesh Venkatesh, and Nojun Kwak. Learning dynamic network using a reuse gate function in semi-supervised video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8405–8414, 2021. 6
- [27] Federico Perazzi, Anna Khoreva, Rodrigo Benenson, Bernt Schiele, and Alexander Sorkine-Hornung. Learning video object segmentation from static images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2663–2672, 2017. 1, 3
- [28] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016. 5
- [29] 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. 5
- [30] Hongje Seong, Junhyuk Hyun, and Euntai Kim. Kernelized memory network for video object segmentation. In *European Conference on Computer Vision*, pages 629–645. Springer, 2020. 1, 2, 3, 6, 7
- [31] Hongje Seong, Seoung Wug Oh, Joon-Young Lee, Seongwon Lee, Suhyeon Lee, and Euntai Kim. Hierarchical memory matching network for video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12889–12898, 2021. 1, 6, 7
- [32] Jianping Shi, Qiong Yan, Li Xu, and Jiaya Jia. Hierarchical image saliency detection on extended cssd. *IEEE transactions on pattern analysis and machine intelligence*, 38(4):717–729, 2015. 6
- [33] Haochen Wang, Xiaolong Jiang, Haibing Ren, Yao Hu, and Song Bai. Swiftnet: Real-time video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1296–1305, 2021. 2, 3, 6, 7
- [34] Lijun Wang, Huchuan Lu, Yifan Wang, Mengyang Feng, Dong Wang, Baocai Yin, and Xiang Ruan. Learning to detect salient objects with image-level supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 136–145, 2017. 6
- [35] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 4
- [36] Huaxin Xiao, Jiashi Feng, Guosheng Lin, Yu Liu, and Maojun Zhang. Monet: Deep motion exploitation for video object segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1140–1148, 2018. 2
- [37] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Shengping Zhang, and Wenxiu Sun. Efficient regional memory network for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1286–1295, 2021. 3, 6, 7
- [38] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018. 5, 6
- [39] Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos K Katsaggelos. Efficient video object segmentation via network modulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6499–6507, 2018. 1, 3
- [40] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by foreground-background integration. In *European Conference on Computer Vision*, pages 332–348. Springer, 2020. 6, 7
- [41] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by multi-scale foreground-background integration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 6, 7
- [42] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4471–4480, 2019. 8
- [43] Yi Zeng, Pingping Zhang, Jianming Zhang, Zhe Lin, and Huchuan Lu. Towards high-resolution salient object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7234–7243, 2019. 6
- [44] Yizhuo Zhang, Zhirong Wu, Houwen Peng, and Stephen Lin. A transductive approach for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6949–6958, 2020. 1, 2