# Look Before You Match: Instance Understanding Matters in Video Object Segmentation

Junke Wang[1,2], Dongdong Chen[3], Zuxuan Wu[1,2], Chong Luo[4], Chuanxin Tang[4],
Xiyang Dai[3], Yucheng Zhao[4], Yujia Xie[3], Lu Yuan[3], Yu-Gang Jiang[1,2]

[1]Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University
[2]Shanghai Collaborative Innovation Center on Intelligent Visual Computing
[3]Microsoft Cloud + AI, [4]Microsoft Research Asia

## Abstract

*Exploring dense matching between the current frame and past frames for long-range context modeling, memory-based methods have demonstrated impressive results in video object segmentation (VOS) recently. Nevertheless, due to the lack of instance understanding ability, the above approaches are oftentimes brittle to large appearance variations or viewpoint changes resulted from the movement of objects and cameras. In this paper, we argue that instance understanding matters in VOS, and integrating it with memory-based matching can enjoy the synergy, which is intuitively sensible from the definition of VOS task, i.e., identifying and segmenting object instances within the video. Towards this goal, we present a two-branch network for VOS, where the query-based instance segmentation (IS) branch delves into the instance details of the current frame and the VOS branch performs spatial-temporal matching with the memory bank. We employ the well-learned object queries from IS branch to inject instance-specific information into the query key, with which the instance-augmented matching is further performed. In addition, we introduce a multi-path fusion block to effectively combine the memory readout with multi-scale features from the instance segmentation decoder, which incorporates high-resolution instance-aware features to produce final segmentation results. Our method achieves state-of-the-art performance on DAVIS 2016/2017 val (92.6% and 87.1%), DAVIS 2017 test-dev (82.8%), and YouTube-VOS 2018/2019 val (86.3% and 86.3%), outperforming alternative methods by clear margins.*

## 1. Introduction

Video object segmentation aims to identify and segment specific objects in a video sequence, which has very broad applications, *e.g.*, interactive video editing and autonomous driving. This work focuses on the semi-supervised setting
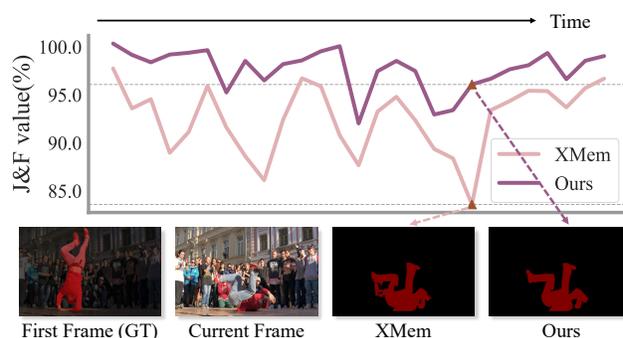


Figure 1. *J&F*-time curve of XMem [13], a state-of-the-art memory-based VOS model and our proposed method. XMem will suffer from a distinct accuracy degradation when the appearance of the target object (*e.g.*, pose of the dancing person) changes dramatically compared to the reference frame. Comparatively, our approach is more robust to this challenging case.

where the annotation of the first frame is given. Starting from Space-Time Memory network (STM) [46], memory-based methods [13–15, 24, 35, 40, 47, 53, 57] have almost dominated this field due to their superior performance and simplicity. STM [46] and its variants [24, 34, 62] typically build a feature memory to store the past frames as well as corresponding masks, and perform dense matching between the query frame and the memory to separate targeted objects from the background.

Despite the prominent success achieved, there exists a non-negligible limitation for the above approaches, *i.e.*, the object deformation and large appearance variations resulted from the motion of camera and objects will inevitably give rise to the risk of false matches [13, 15, 46], thus making them struggle to generate accurate masks. We visualize the *J&F*-time curve of XMem [13], a state-of-the-art memory-based VOS model, on a representative video from DAVIS 2017 in Figure 1. It can be seen that, when the target object undergoes a distinct pose change compared to the first
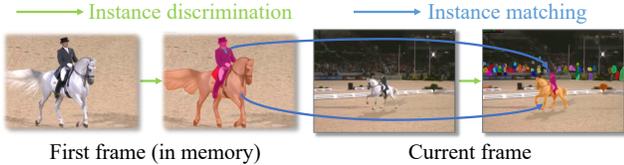
Figure 2. **A conceptual introduction on how humans address the VOS task**. For the current frame of a video stream, humans first distinguish between different instances and then match them with the target object(s) in memory.

reference frame, XMem [13] misidentifies another person wearing the same color as the foreground and suffers from drastic performance degradation.

In contrast, humans are capable of avoiding such mistakes and achieving consistently accurate matching. This gap motivates us to reflect on how we humans resolve the VOS task. Intuitively, given the current frame in a video sequence, humans typically first distinguish between different instances within it by identifying which instance each pixel belongs to. After that, the instance matching with the target object(s) in memory is conducted to obtain the final results (see Figure 2 for a conceptual illustration). In fact, this intuition is also consistent with the definition of VOS itself, *i.e.*, identify (matching) and segmenting objects (instance understanding). Moreover, in the absence of instance understanding, it is theoretically difficult to generate accurate predictions for regions that are invisible in reference frame by pure matching.

Inspired by this, we argue that instance understanding is critical to the video object segmentation, which could be incorporated with memory-based matching to enjoy the synergy. More specifically, we aim to derive instance-discriminative features that are able to ***distinguish different instances***. Equipped with these features, we then perform ***semantic matching*** with the memory bank to effectively associate the target object(s) with specific instance(s) in the current frame.

In this spirit, we present a two-branch network, ISVOS, for semi-supervised VOS, which contains an instance segmentation (IS) branch to delve into the instance details for the current frame and a video object segmentation branch that resorts to an external memory for spatial-temporal matching. The IS branch is built upon a query-based instance segmentation model [10] and supervised with instance masks to learn instance-specific representations. Note that ISVOS is a generic framework and IS branch can be easily replaced with more advanced instance understanding models. The video object segmentation (VOS) branch, on the other hand, maintains a memory bank to store the features of past frames and their predictions. We compare the query key of current frame and memory key[1] from mem-

ory bank for affinity calculation following [13, 15, 46, 53]. Motivated by recent approaches that use learnable queries serving as region proposal networks to identify instances in images [10, 11, 19, 68], we employ object queries from the IS branch to inject instance-specific information into our query key, with which the instance-augmented matching is performed. After that, the readout features are produced by aggregating the memory value with the affinity matrix. Moreover, in order to make use of the fine-grained instance details reserved in high-resolution instance-aware features, we further combine the multi-scale features from instance segmentation decoder with the memory readout through a carefully designed multi-path fusion block to finally generate the segmentation masks.

We conduct experiments on the standard DAVIS [49,50] and YouTube-VOS [64] benchmarks. The results demonstrate that our ISVOS can achieve state-of-the-art performance on both single-object (*i.e.*, 92.6% in terms of $J\&F$ on DAVIS 2016 validation split) and multi-object benchmarks (*i.e.*, 87.1% and 82.8% on DAVIS 2017 validation and test-dev split, 86.3% and 86.3% on YouTube-VOS 2018 & 2019 validation split) without post-processing.

## 2. Related Work

**Propagation-based VOS.** Propagation-based VOS methods [16, 18, 31, 45, 55, 63, 66, 67] take advantage of temporal correlations between adjacent frames to iteratively propagate the segmentation masks from the previous frame to the current frame. Early approaches [3, 25, 48] typically follow an online learning manner by finetuning models at test-time, which therefore suffer from limited inference efficiency. To mitigate this issue, the following studies shift attention to offline learning by utilizing optical flow [16,55,66] as guidance to deliver the temporal information smoothly. Despite the promising results achieved, these methods are oftentimes vulnerable to the error accumulation brought by occlusion or drifting.

**Matching-based VOS.** In order to model the spatial-temporal context over longer distances, matching-based models [14, 15, 26, 46] typically calculate the correspondence between the current frame and the reference frame [8, 26, 56, 69] and even a feature memory [13, 24, 35, 40, 52, 53, 62] to identify target objects. In addition, several studies focus on designing novel memory construction [34, 35] or matching [57] strategies to improve the inference efficiency of VOS models. However, the widely adopted dense matching between reference features will inevitably fail on the objects with significant appearance variations or viewpoint changes. In this paper, we propose to integrate in-

---

[1]In this paper, we follow previous work [15,46] of which are compared

with query key to denote the key features of current frame and memory bank as query key and memory key, respectively, so as to perform instance-augmented matching.
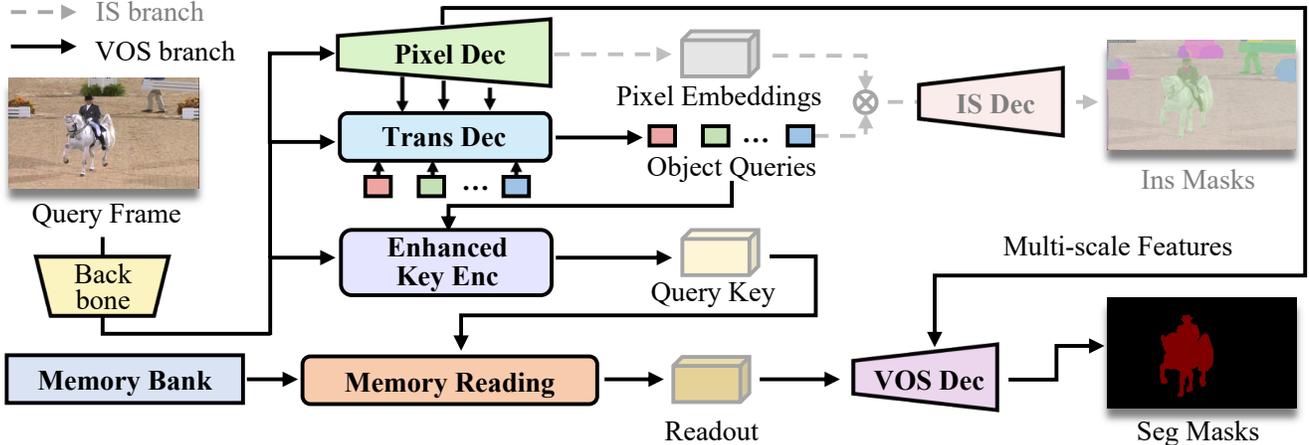
Figure 3. Overview of the proposed method, which consists of an **instance segmentation branch** and a **video object segmentation branch**. We jointly train both branches on instance segmentation and video object segmentation tasks, respectively. The IS branch parts denoted as dotted gray line will be skipped during inference, meaning that our method does not explicitly use the output instance masks.

stance understanding into VOS, which is neglected in existing matching-based methods.

**Instance Segmentation.** Built upon powerful object detectors, two-stage instance segmentation models [2, 7, 22, 27, 33, 38] predict bounding boxes first and then extract the instance mask in each region of interest (ROI). These methods require complicated procedures, *e.g.*, region proposal generation [51] and ROIAlign [22], which motivates the following work to develop one-stage box-free models [1, 9, 20, 29, 37, 43, 44, 59, 60]. Recently, the success of DETR [4] inspires a series of query-based models [10, 11, 19] to reformulate the instance segmentation from a novel "set prediction" perspective, which achieve state-of-the-art performance on standard benchmarks [36]. In this work, we build our VOS model with an auxiliary query-based instance segmentation model [10] to make the intermediate feature instance-aware rather than explicitly use the output segmentation mask. It enables us to perform instance-augmented matching between the instance-aware features of the current frame and memory. Note that our approach is generic and can be employed to improve various memory-based VOS methods.

# 3. Method

Our goal is to integrate the instance understanding for improved memory-based video object segmentation. To this end, we propose ISVOS, a two-branch network where the instance segmentation (IS) branch learns instance-specific information, and the video object segmentation (VOS) branch performs instance-augmented matching through memory reading. Such a design also shares the similar spirit with a prior study [21] which implies that two layers of human cortex responsible for the object recognition and track-

ing separately.

More formally, given a video sequence $\mathcal{V} = [X_1, X_2, ..., X_T]$ and the annotation mask of the first frame, we process the frames sequentially and maintain a memory bank to store the past frames and their predictions following [15, 46]. For the current frame $X_t \in \mathcal{R}^{3 \times H \times W}$, we first extract $F_{res4}$ from ResNet [23] as our backbone features, which is shared by a pixel decoder to generate per-pixel embeddings and a Transformer decoder to inject localized features to learnable object queries in IS branch. While in the VOS branch, we apply an Enhanced Key Encoder to project backbone feature $F_{res4}$ to query key, which are compared with the memory key to perform semantic matching. Finally, the memory readout as well as multi-scale features from both backbone and pixel decoder are input to the VOS decoder to produce the final mask prediction. The architecture of ISVOS is illustrated in Figure 3. Below, we first introduce IS and VOS branch in Sec. 3.1 and Sec. 3.2, respectively, and then elaborate how to obtain the final mask prediction in Sec 3.3.

## 3.1. Instance Segmentation Branch

As is described above, existing memory-based VOS models typically perform dense matching between the features of current frame and memory bank without mining the instance information, which therefore suffers from false matches when distinct object deformation or appearance changes happen. To address this issue, we explore to acquire the instance understanding capability from an instance segmentation (IS) branch, which is built upon an auxiliary query-based instance segmentation model [10]. Specifically, our IS branch consists of the following components:

**Pixel Decoder** takes $F_{res4}$ as input and generates per-pixel embeddings $F_{pixel} \in \mathbb{R}^{C_\varepsilon \times H/4 \times W/4}$ with alternate convo-

lutional and upsampling layers, where $C_\varepsilon$ is the embedding dimension. In addition, we also input the feature pyramid $\{P_i\}_{i=0}^2$ produced by the pixel decoder with a resolution 1/32, 1/16 and 1/8 of the original image into both the transformer decoder and the VOS decoder, so as to fully take advantage of high-resolution instance features. For each resolution, we add both a sinusoidal positional embedding and a learnable scale-level embedding following [10, 74].

**Transformer Decoder** gathers the local information in features obtained by pixel decoder to a set of learnable object queries $q_{ins} \in \mathcal{R}^{N \times C_d}$ through masked attention [10], where $N$ is a pre-defined hyper-parameter to indicate the number of object queries (which is set to 100, empirically), and $C_d$ is the query feature dimension. The masked attention can be formulated as:

$$q_l = \text{softmax}(\mathcal{M}_{l-1} + q_l k_l^{\text{T}})v_l + q_{l-1}, \qquad (1)$$

where $k_l$ and $v_l$ are the key and value embeddings projected from one resolution of $\{P_i\}_{i=0}^2$ (with $i$ corresponds to $l$), respectively. We add an auxiliary loss [10] to every intermediate Transformer decoder layer and $\mathcal{M}_{l-1}$ is the binarized (threshold = 0.5) mask prediction from the previous $(l-1)_{th}$ layer. Note that $q_0$ is initialized with $q_{ins}$ and finally updated to $\tilde{q}_{ins}$ as the final object queries.

### 3.2. Video Object Segmentation Branch

With the instance-specific information from the IS branch, our VOS branch further performs instance-augmented matching between the current frame and the memory bank in the feature space, so as to leverage long-range context information for mask generation.

**Enhanced Key Encoder** takes the updated object queries $\tilde{q}_{ins}$ and the backbone feature $F_{res4}$ as inputs to generate the query key of current frame with $C_k$-dimension. Specifically, we follow [15] to first apply a $3 \times 3$ convolutional layer on top of $F_{res4}$ to obtain $Q_g \in \mathcal{R}^{C_h \times H/16 \times W/16}$, $C_h$ denotes the hidden dimension and we set it to the query feature dimension. Next, we aggregate the image features in $Q_g$ to $\tilde{q}_{ins}$ through a deformable attention layer [74]:

$$\tilde{q}_{vos} = \text{DeformAttn}(\tilde{q}_{ins}, p, Q_g), \qquad (2)$$

where $p$ is a set of 2-d reference points.

After that, we inject instance information in $\tilde{q}_{vos}$ to $Q_g$ reversely through a dot product (flattening operation is omitted for brevity), which is then activated by a sigmoid function and concatenated with $Q_g$ to get $Q_{cat}$. Finally, we apply a convolutional projection head on $Q_{cat}$ and further flatten it to the instance-aware query key $Q \in \mathcal{R}^{C_k \times H_m W_m}$, where $H_m = H/16$ and $W_m = W/16$.

**Memory Reading** first retrieves the memory key $K \in \mathcal{R}^{C_K \times T H_m W_m}$ and memory value $V \in \mathcal{R}^{C_v \times T H_m W_m}$ from memory bank, where $T$ is the current memory size,

and $C_v$ denotes the value feature dimension. Then the similarity between $K$ and the above query key $Q$ is measured by calculating the affinity matrix:

$$A_{i,j} = \frac{\exp(\text{d}(K_i, Q_j))}{\sum_i(\exp(\text{d}(K_i, Q_j)))}, \qquad (3)$$

where the subscript indexes the spatial location of query and memory key, and the distance function here we use is L2 distance following [13, 15]. Note that we normalize the distance value by $\sqrt{C_k}$ as in [15, 46].

With the affinity matrix $A$, the memory value $V$ could be aggregated through a weighted summation to obtain the readout features $F_{mem} \in \mathcal{R}^{C_v \times H_m W_m}$. Finally, we pass $F_{mem}$ to the VOS decoder for mask generation, which will be clarified later.

**Memory Update** is executed once the prediction of the current frame is generated during training and at a fixed interval during inference, which stores the memory key and memory value of current frame into the memory bank. We follow [13, 15] to simply share the key between query and memory, *i.e.*, the query key will be saved in the memory bank as a memory key if the current frame should be "memorized". While for the memory value, we first input the predicted mask to a lightweight backbone (ResNet18 [23] is adopted in this paper), the last layer feature of which is concatenated with $F_{res4}$ to obtain $\check{V}_{cur}$ following [13, 15]. Next, we further input $\check{V}_{cur}$ to two ResBlocks and a CBAM block [61] sequentially to get the memory value $V_{cur} \in \mathcal{R}^{C_v \times H_m W_m}$. In this paper, we describe the forward process for a single target object for readability. In the case of multi-object segmentation, an extra dimension is needed for $V$ to indicate the number of objects [15].

### 3.3. Mask Prediction

On top of the IS branch and VOS branch, we apply an auxiliary instance segmentation decoder and a video object segmentation decoder to generate the instance mask and video object mask predictions, respectively. Note that *the auxiliary instance segmentation decoder along with the pixel embedding will only be used during training and discarded during inference.*

**Instance Segmentation Decoder** inputs the updated object queries $\tilde{q}_{ins}$ to a linear classifier and a softmax activation function successively to yield category probabilities. Besides, a Multi-Layer Perceptron (MLP) with 2 hidden layers transforms $\tilde{q}_{ins}$ to the corresponding mask embeddings. Finally, we obtain each binary mask prediction $\hat{M}_{ins}$ via a dot product between the mask embedding and per-pixel embeddings $F_{pixel}$.

**Video Object Segmentation Decoder** fuses the memory readout $F_{mem}$, multi-scale features from backbone

| Method | w/ BL30K | DAVIS16 validation | | | DAVIS17 validation | | | YT2018 validation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{G}$ | $\mathcal{J}_s$ | $\mathcal{F}_s$ | $\mathcal{J}_u$ | $\mathcal{F}_u$ |
| STM [46] | ✗ | 89.3 | 88.7 | 89.9 | 81.8 | 79.2 | 84.3 | 79.4 | 79.7 | 84.2 | 72.8 | 80.9 |
| HMMN [54] | ✗ | 90.8 | 89.6 | 92.0 | 84.7 | 81.9 | 87.5 | 82.6 | 82.1 | 87.0 | 76.8 | 84.6 |
| RPCM [65] | ✗ | 90.6 | 87.1 | 91.1 | 83.7 | 81.3 | 86.0 | 84.0 | 83.1 | 87.7 | 78.5 | 86.7 |
| STCN [15] | ✗ | 91.6 | 90.8 | 92.5 | 85.4 | 82.2 | 88.6 | 83.0 | 81.9 | 86.5 | 77.9 | 85.7 |
| AOT [70] | ✗ | 91.1 | 90.1 | 92.1 | 84.9 | 82.3 | 87.5 | 85.5 | 84.5 | 89.5 | 79.6 | 88.2 |
| RDE [30] | ✗ | 91.1 | 89.7 | 92.5 | 84.2 | 80.8 | 87.5 | - | - | - | - | - |
| XMem [13] | ✗ | 91.5 | 90.4 | 92.7 | 86.2 | 82.9 | 89.5 | 85.7 | 84.6 | 89.3 | 80.2 | 88.7 |
| DeAOT [72] | ✗ | 92.3 | 90.5 | 94.0 | 85.2 | 82.2 | 88.2 | 86.0 | 84.9 | 89.9 | 80.4 | 88.7 |
| Ours | ✗ | **92.6** | **91.5** | **93.7** | **87.1** | **83.7** | **90.5** | **86.3** | **85.5** | **90.2** | **80.5** | **88.8** |
| MiVOS [14] | ✓ | 91.0 | 89.6 | 92.4 | 84.5 | 81.7 | 87.4 | 82.6 | 81.1 | 85.6 | 77.7 | 86.2 |
| STCN [15] | ✓ | 91.7 | 90.4 | 93.0 | 85.3 | 82.0 | 88.6 | 84.3 | 83.2 | 87.9 | 79.0 | 87.3 |
| RDE [30] | ✓ | 91.6 | 90.0 | 93.2 | 86.1 | 82.1 | 90.0 | - | - | - | - | - |
| XMem [13] | ✓ | 92.0 | 90.7 | 93.2 | 87.7 | 84.0 | 91.4 | 86.1 | 85.1 | 89.8 | 80.3 | **89.2** |
| Ours | ✓ | **92.8** | **91.8** | **93.8** | **88.2** | **84.5** | **91.9** | **86.7** | **86.1** | **90.8** | **81.0** | 89.0 |

Table 1. Quantitative comparisons on the DAVIS 2016 val, DAVIS 2017 val, and YouTube-VOS 2018 val split.
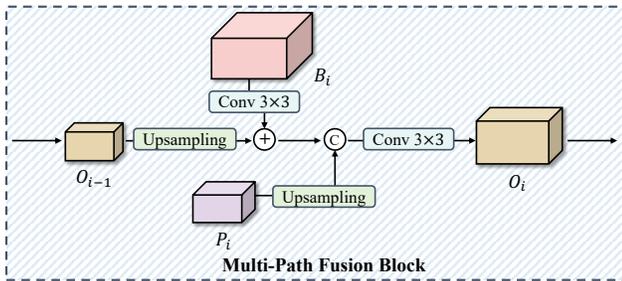


Figure 4. Illustration of the Multi-Path Fusion (MPF) block.

$\{B_i\}_{i=0}^2$[2] and pixel decoder $\{P_i\}_{i=0}^2$ with a multi-path fusion (MPF) block to make use of the fine-grained details reserved in high-resolution instance-aware features. The MPF block could be formulated as:

$$O_i = \text{MPF}(O_{i-1}, B_i, P_i), \qquad (4)$$

$O_{i-1}$ is output by the previous MPF block, which is initialized with $F_{mem}$. Specifically, we first input $B_i$ and $P_i$ to $3 \times 3$ convolutional layers to align their feature dimensions with $O_{i-1}$ to obtain $\tilde{B}_i$ and $\tilde{P}_i$, separately. Next, we concatenate the sum of upsampled $O_{i-1}$ and $\tilde{B}_i$ with the upsampled $\tilde{P}_i$, the result of which is finally input to a residual block to get $O_i$. The detailed architecture of MPF block is illustrated in Figure 4. Note that the final layer of the decoder produces a mask with stride = 4, and we bilinearly upsample it to the original resolution.

---

[2]We follow the previous work [13, 15, 46] to adopt the features with stride = 4, 8, and 32.

# 4. Experiments

## 4.1. Implementation Details

**Training.** The instance segmentation (IS) branch and video object segmentation (VOS) branch are jointly trained with different supervisory signals. We train the IS branch on a standard instance segmentation dataset COCO [36] with a resolution of 384 (compatible with VOS branch). While for the VOS branch, we follow [13, 15, 35, 46, 52] to firstly pretrain our network on deformed static images [12, 32, 54, 58, 73], and then perform the main training on YouTube-VOS [64] and DAVIS [50]. Note that we also pretrain on BL30K [6, 14, 17] optionally to further boost the performance following [13, 15], and the models pretrained on additional data are denoted with an asterisk (∗).

The IS branch is supervised with a combination of mask loss and classification loss, where the mask loss consists of weighted binary cross-entropy loss and dice loss [42]. The VOS branch is supervised with bootstrapped cross entropy loss and dice loss following [70]. The static image pretraining lasts 150K iterations with a batch size of 56 and a learning rate of 4e-5. While the main training lasts 110K iterations with a batch size 16 and a learning rate 2e-5. The complete model is optimized with AdamW [28, 39] with a weight decay of 0.05. We load the COCO-pretrained Mask2Former [10] to initialize our instance segmentation branch, and use $0.1\times$ learning rate for these parameters. The overall learning rate is decayed by a factor of 10 after the first 80K iterations.

**Inference.** We follow previous work [13, 15, 46, 52] to memorize every 5th frame during inference. Specially, we

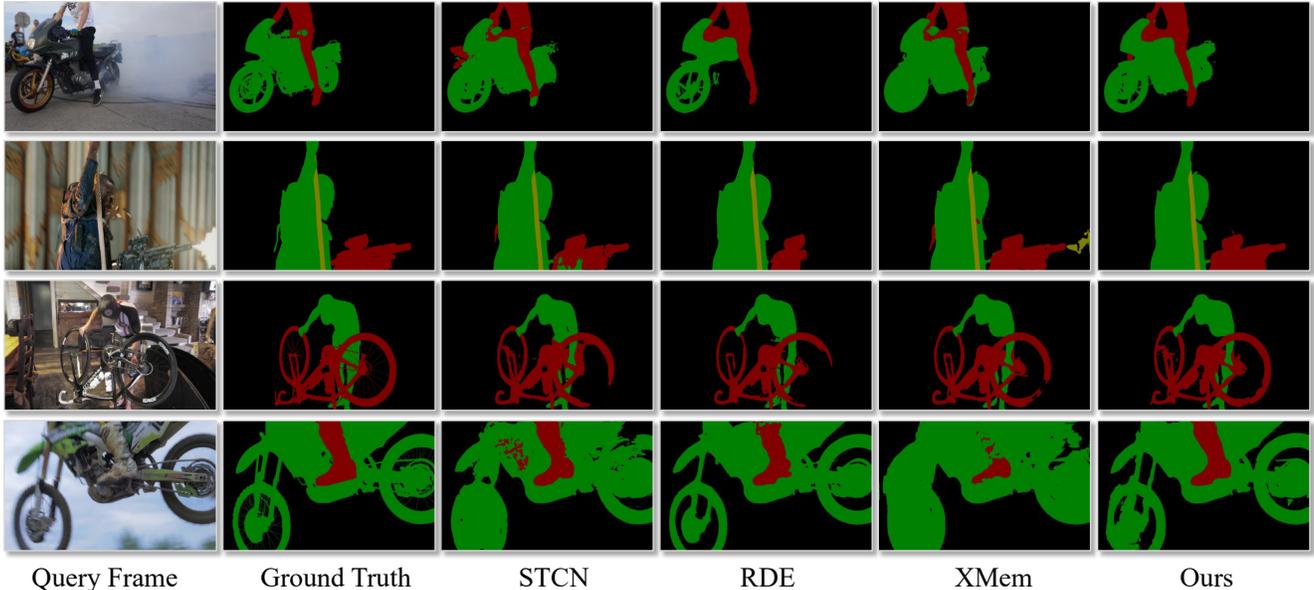| Query Frame | Ground Truth | STCN | RDE | XMem | Ours |

Figure 5. Qualitative comparisons between ISVOS and several state-of-the-art memory-based VOS models, including RDE [30], STCN [15], and XMem [13].

implement the memory as a first-in-first-out (FIFO) queue, and restrict the maximum memory size to 16 to improve the inference speed. Note that the first frame and its mask are always reserved to provide accurate reference information. We adopt Top-K filter [13–15] for memory-reading augmentation, with K set to 20.

**Evaluation datasets and metrics.** We evaluate the performance of ISVOS on standard VOS datasets DAVIS [49, 50] and YouTube-VOS [64]. DAVIS 2016 [49] is a single-object VOS benchmark and DAVIS 2017 [50] extends it to a multi-object version. YouTube-VOS [64] is the large-scale benchmark for multi-object VOS, which also includes unseen categories in the validation set to measure the generalization ability. We report the results on 474 and 507 validation videos in the 2018 and 2019 versions (denoted as "YT2018" and "YT2019" in the following Tables respectively). We use mean Jaccard $\mathcal{J}$ index and mean boundary $\mathcal{F}$ score, along with mean $\mathcal{J}\&\mathcal{F}$ to evaluate segmentation accuracy. Note that for YouTube-VOS, we report the results on both seen and unseen categories, along with the averaged overall score $\mathcal{G}$.

### 4.2. Comparison with State-of-the-art Methods

The comparison results between ISVOS and existing state-of-the-art VOS models on DAVIS 2016 validation, DAVIS 2017 validation, and YouTube-VOS 2018 validation are listed in Table 1. We can see that without incorporating BL30K as addition training data, our method achieves top-ranked performance on both single-object and multi-object VOS benchmarks, i.e., 92.6%, 87.1%, 86.3% in terms of

| Method | D17 test-dev | | | YT2019 validation | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{G}$ | $\mathcal{J}_s$ | $\mathcal{F}_s$ | $\mathcal{J}_u$ | $\mathcal{F}_u$ |
| HMMN [53] | 78.6 | 74.7 | 82.5 | 82.5 | 81.7 | 86.1 | 77.3 | 85.0 |
| STCN [15] | 76.1 | 73.1 | 80.0 | 82.7 | 81.1 | 85.4 | 78.2 | 85.9 |
| RPCM [65] | 79.2 | 75.8 | 82.6 | 83.9 | 82.6 | 86.9 | 79.1 | 87.1 |
| AOT [70] | 79.6 | 75.9 | 83.3 | 85.3 | 83.9 | 88.8 | 79.9 | 88.5 |
| RDE [30] | 77.4 | 73.6 | 81.2 | 81.9 | 81.1 | 85.5 | 76.2 | 84.8 |
| XMem [13] | 81.0 | 77.4 | 84.5 | 85.5 | 84.3 | 88.6 | 80.3 | 88.6 |
| DeAOT [72] | 80.7 | 76.9 | 84.5 | 85.9 | 84.6 | 89.4 | 80.8 | 88.9 |
| Ours | **82.8** | **79.3** | **86.2** | **86.1** | **85.2** | **89.7** | **80.7** | **88.9** |
| MiVOS* [14] | 78.6 | 74.9 | 82.2 | 82.4 | 80.6 | 84.7 | 78.1 | 86.4 |
| STCN* [15] | 77.8 | 74.3 | 81.3 | 84.2 | 82.6 | 87.0 | 79.4 | 87.7 |
| RDE [30] | 78.9 | 74.9 | 82.9 | 83.3 | 81.9 | 86.3 | 78.0 | 86.9 |
| XMem* [13] | 81.2 | 77.6 | 84.7 | 85.8 | 84.8 | 89.2 | 80.3 | 88.8 |
| Ours* | **84.0** | **80.1** | **87.8** | **86.3** | **85.2** | **89.7** | **81.0** | **89.1** |

Table 2. Results on DAVIS 2017 (D17) test-dev and YouTube-VOS 2019 validation. * denotes BL30K is adopted for pretraining.

$\mathcal{J}\&\mathcal{F}$ on DAVIS 2016 & 2017, YouTube-VOS 2018 validation split, respectively, even surpassing existing methods that are pretrained on BL30K. Adopting BL30K as additional training data can further boost the performance of ISVOS. We also report the results on DAVIS 2017 test-dev and YouTube-VOS 2019 validation split in Table 2, and ISVOS also outperforms all the baseline methods. Even though our method adopts a simpler memory mechanism than existing methods like [13,53,65], we still achieve superior performance. This highlights that introducing instance understanding to conduct instance-augmented matching is

super helpful and clearly outperforms the vanilla semantic matching used in the existing methods [13, 15, 46].

We further visualize the segmentation results of ISVOS on some representative challenging cases with dramatic movements (*e.g.*, shooting and motor cross-jump), and compare with state-of-the-art memory-based VOS models including STCN [15], RDE [30], and XMem [13] in Figure 5. We can see that RDE [30] and STCN [15] struggle with occlusions incurred by smoke and confusing objects, respectively. XMem [13] produces more competitive results, which however fails to generate sharp boundaries for the motorcycle rims. Our method, by contrast, generates more accurate and clear masks on these challenging cases. This suggests that the instance-aware representations learned from the instance segmentation branch could facilitate our model to derive instance-discriminative features.

### 4.3. Discussion

**Impact of query enhancement and MPF block.** The query enhancement (QE) (Equation 2) is used to inject instance-specific into query key, while the multi-path fusion (MPF) block (Equation 4) is designed to incorporate high-resolution instance-aware features for fine-grained detail prediction. To evaluate their effectiveness, we remove QE and MPF separately from ISVOS and evaluate the segmentation performance on DAVIS 2017 validation (DAVIS17 val) and YouTube-VOS 2018 validation (YT2018 val) split.

| Method | DAVIS17 val | | | YT2018 val | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{G}$ | $\mathcal{J}_s$ | $\mathcal{F}_s$ | $\mathcal{J}_u$ | $\mathcal{F}_u$ |
| w/o QE & MPF | 85.3 | 82.0 | 88.6 | 83.0 | 84.0 | 75.7 | 88.5 | 83.8 |
| w/o QE | 85.7 | 82.4 | 88.9 | 84.4 | 85.1 | 77.4 | 89.8 | 85.5 |
| w/o MPF | 86.2 | 83.0 | 89.5 | 85.6 | 85.0 | 90.4 | 79.4 | 87.5 |
| Ours | **87.1** | **83.7** | **90.5** | **86.3** | **85.5** | **90.2** | **80.5** | **88.8** |

Table 3. Results on DAVIS 2017 validation and YouTube-VOS validation split w/ and w/o query enhancement (QE) and multi-path fusion (MPF) block.

The quantitative results are listed in Table 3. We can observe that without QE, the $\mathcal{J}\&\mathcal{F}$ value decreases by 1.4% on DAVIS17 val and 1.9% on YT2018 val, while without MPF, the $\mathcal{J}\&\mathcal{F}$ value decreases by 0.9% and 0.7%, respectively. The performance degradation validates that the use of the above components both effectively improves the performance of our model.

**Impact of weight initialization and joint training for the IS branch.** The IS branch is built upon a instance segmentation model [10] to acquire instance-specific information. In our implementation, we load the weights from Mask2Former [10] pretrained on COCO [36] and perform joint training on both IS task and VOS task to prevent the catastrophic forgetting. To study the effect of weight ini-

tialization and joint training, we conduct experiments under different settings and compare the results in Table 4.

| LW | JT | DAVIS17 val | | | YT2018 val | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{G}$ | $\mathcal{J}_s$ | $\mathcal{F}_s$ | $\mathcal{J}_u$ | $\mathcal{F}_u$ |
| ✗ | ✗ | 78.6 | 75.6 | 81.5 | 78.6 | 79.2 | 83.5 | 72.3 | 79.4 |
| ✗ | ✓ | 80.0 | 76.9 | 83.1 | 80.6 | 80.1 | 84.5 | 74.5 | 83.2 |
| ✓ | ✗ | 82.0 | 79.2 | 84.4 | 81.3 | 80.7 | 85.4 | 75.5 | 83.6 |
| ✓ | ✓ | **87.1** | **83.7** | **90.5** | **86.3** | **85.5** | **90.2** | **80.5** | **88.8** |

Table 4. Results on DAVIS 2017 validation and YouTube-VOS validation split w/ and w/o loading the weights from pretrained Mask2Former [10] (**LW**) and joint training (**JT**).
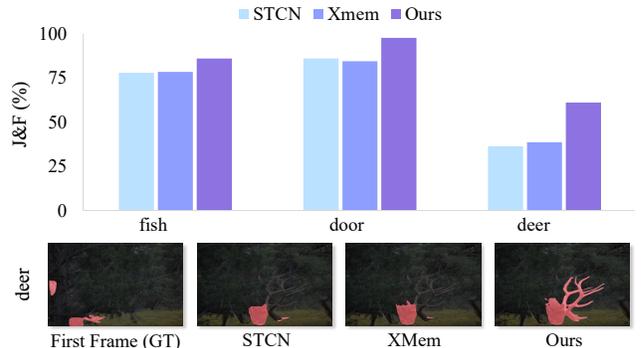


Figure 6. $\mathcal{J}\&\mathcal{F}$ metric of STCN [15], XMem [13], and ISVOS on three videos whose categories do not appear in the COCO dataset.

The drastic performance drop in the first row indicates that instance-aware representations are critical to the generation of accurate masks. Performing joint training from scratch brings about minor improvements, but it is difficult for the VOS branch to learn useful instance information from IS branch in the beginning. Initializing the weights from Mask2Former [10] improves the performance more significantly, which however, would gradually lose the instance segmentation capability without joint training. In contrast, the combined weight initialization from Mask2Former and joint training achieves the best results. We would also like to point out while we use Mask2Former for initialization, the IS branch can be easily replaced with any query-based instance segmentation model.

In addition, considering that VOS is essentially a category-agnostic task but the IS branch is trained on a close set, we further show the performance of ISVOS on objects that do not appear in COCO, *e.g.*, fish, door, and deer[3], and compare with STCN and XMem in Figure 6. The quantitative comparison is also displayed. We can see that our method still performs well on these objects and generates more accurate masks than existing methods. This indicates

---

[3]These objects correspond to the f78b3f5f34, 4d6cd98941, and f6ed698261 video in YouTube-VOS 2018 val.

that joint training allows our method to develop generalizable instance differentiation capability even if the IS branch is trained on a close-set instance segmentation dataset.

**Trade-off between memory size and segmentation performance.** As mentioned in Sec. 4.1, we implement the memory bank as a first-in-first-out (FIFO) queue with a maximum size. To further investigate the behavior of ISVOS, we dynamically adjust the maximum memory size and observe the trend of performance variation (*i.e.*, $\mathcal{J}\&\mathcal{F}$) on DAVIS 2016 & 2017 validation split. We also re-implement the memory bank of several existing memory-based models (including STM [46], STCN [15], and XMem [13]) as FIFO queues, and compare their results with ISVOS in Figure 7.
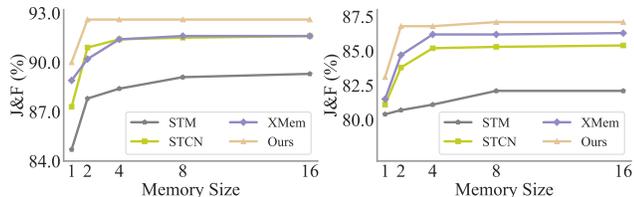


Figure 7. Tradeoffs between the maximum memory size and $\mathcal{J}\&\mathcal{F}$ on DAVIS 2016 (left) and 2017 (right) validation split.

We can see that increasing the memory size always improves the segmentation performance for all these methods, since more contextual information is used. It is noteworthy that, our method achieves competitive results by relying on a smaller memory size, *e.g.*, when the memory size is set to 2, the $\mathcal{J}\&\mathcal{F}$ value is only 0.3 away from the highest point for ISVOS, but 1.6 for XMem on DAVIS 2017 validation split. This demonstrates the superiority of instance-augmented matching compared with vanilla semantic matching.

**Results with different training data.** In the main experiment, we follow previous methods [13–15, 46] to first pretrain our model on static images (and BL30K optionally) for fair comparisons. To study the effects of pretraining on the final segmentation results, we additionally conduct experiments to train ISVOS on DAVIS 2017 [50] only, YouTube-VOS 2019 [64] only, and a mix of both. The comparison with existing models are shown in Table 5. We can see that ISVOS achieves competitive results even without incorporating static images and BL30K for pretraining, outperforming all the baseline models by a large margin. When gradually increasing the scale of the training data, the performance of our method can be further boosted.

**Multi-scale Inference.** Multi-scale evaluation is a commonly used trick in segmentation tasks [5, 13, 15] to boost the performance by merging the results of inputs under different data augmentations. Here we follow XMem [13] to apply image scaling and vertical mirroring and simply aver-

| Method | DAVIS17 val | | | YT2018 val | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{G}$ | $\mathcal{J}_s$ | $\mathcal{F}_s$ | $\mathcal{J}_u$ | $\mathcal{F}_u$ |
| SST‡ [18] | 82.5 | 79.9 | 85.1 | 81.7 | 81.2 | - | 76.0 | - |
| CFBI+‡ [71] | 82.9 | 80.1 | 85.7 | 82.0 | 81.2 | 86.0 | 76.2 | 84.6 |
| JOINT‡ [41] | 83.5 | 80.8 | 86.2 | 83.1 | 81.5 | 85.9 | 78.7 | 86.5 |
| XMem‡ | 84.5 | - | - | 84.3 | - | - | - | - |
| Ours‡ | 85.2 | 82.1 | 88.3 | 84.7 | 84.5 | 89.1 | 78.2 | 87.0 |
| D only | 77.5 | 75.6 | 79.4 | - | - | - | - | - |
| Y only | - | - | - | 84.9 | 84.0 | 88.8 | 78.8 | 88.0 |
| S + D + Y | 87.1 | 83.7 | 90.5 | 86.3 | 85.5 | 90.2 | 80.5 | 88.8 |
| S + D + B + Y | 88.2 | 84.5 | 91.9 | 86.7 | 86.1 | 90.8 | 81.0 | 89.0 |

Table 5. Results on DAVIS 2017 validation and YouTube-VOS validation split with different training data. D: DAVIS 2017, Y: YouTube 2019, S: static images, B: BL30K. ‡ denotes pretraining on the combined DAVIS and YouTube-VOS data (*i.e.*, D + Y).

age the output probabilities to obtain the final masks.

| Method | MS | D16 val | | | D17 val | | |
|---|---|---|---|---|---|---|---|
| | | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
| CFBI [69] | ✓ | 90.7 | 89.6 | 91.7 | 83.3 | 80.5 | 86.0 |
| XMem [13] | ✓ | 92.7 | 92.0 | 93.5 | 88.2 | 85.4 | 91.0 |
| Ours | ✗ | 92.6 | 91.5 | 93.7 | 87.1 | 83.7 | 90.5 |
| Ours | ✓ | 92.9 | 92.2 | 93.6 | 88.6 | 85.8 | 91.4 |
| Ours* | ✗ | 92.8 | 91.8 | 93.8 | 88.2 | 84.5 | 91.9 |
| Ours* | ✓ | 93.4 | 92.5 | 94.2 | 89.8 | 86.7 | 93.0 |

Table 6. Results on DAVIS 2017 validation and YouTube-VOS validation split with different training data. D: DAVIS 2017, Y: YouTube 2019, S: static images, B: BL30K. ‡ denotes pretraining on the combined DAVIS and YouTube-VOS data.

The results in Table 8 imply that multi-scale inference improves the performance of ISVOS by 0.3% and 1.7% in terms of $\mathcal{J}\&\mathcal{F}$ on DAVIS 2016 / 2017 validation split, and ISVOS still outperforms existing methods.

**Results on Long video datasets.** In order to further evaluate the long-term performance of ISVOS, we additionally test our method on the Long-time Video dataset [35], which contains three videos with more than 7,000 frames in total for validation. Considering the video duration is longer and the target object(s) will undergo distinct appearance deformation or scale variations, we set the maximum memory size to 64 during inference. The comparison results are shown in Table 9.

We can observe that ISVOS again achieves the best segmentation results measured in different metrics. It is worth mentioning that ISVOS beats the methods specifically designed for long videos, *e.g.*, AFB-URR [35] and

| Method | Long-time Video | | |
|--------|:----:|:----:|:----:|
| | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
| RMNet [62] | 59.8 | 59.7 | 60.0 |
| JOINT [41] | 67.1 | 64.5 | 69.6 |
| STM [46] | 80.6 | 79.9 | 81.3 |
| HMMN [54] | 81.5 | 79.9 | 83.0 |
| STCN [15] | 87.3 | 85.4 | 89.2 |
| AOT [70] | 84.3 | 83.2 | 85.4 |
| AFB-URR [35] | 83.7 | 82.9 | 84.5 |
| XMem [13] | 89.8 | 88.0 | 91.6 |
| Ours | **90.0** | **88.3** | **91.7** |

Table 7. Results on the Long-time Video dataset [35].

XMem [13]. We believe the performance gain is resulted from taking advantage of the instance information in the query frame to facilitate the semantic matching.

**Visualization of Readout Features.** We visualize the readout features (*i.e.*, $F_{mem}$ in Sec. 3.2) of several memory-based VOS models and ISVOS in Figure 8 to further compare the vanilla semantic matching without instance understanding and instance-augmented matching. The high resolution feature $P_2$ from the pixel-decoder (Sec. 3.1) is also displayed. We can see that with the enhanced query key, the instance information in our readout features is more clear and distinguishable. In addition, the abundant details in high-resolution instance-aware features also help ISVOS to produce sharp boundaries.



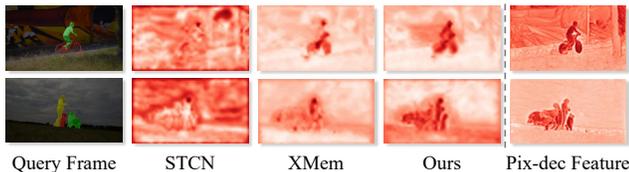Query Frame    STCN    XMem    Ours    Pix-dec Feature

Figure 8. Visualization of the query frame, memory readout features of STCN [15], XMem [13], our method, and the instance-aware features of the highest resolution $P_2$ from the pixel-decoder.

## 5. Conclusion

This paper proposes to incorporate instance understanding into memory-based matching for improved video object segmentation. To achieve this goal, a two-branch network ISVOS is introduced, where the instance segmentation (IS) branch derives instance-aware representations of current frame and the video object segmentation (VOS) branch maintains a memory bank for spatial-temporal matching. We enhance the query key with the well-learned object queries from IS branch to inject the instance-specific information, with which the instance-augmented matching with

memory bank is performed. Furthermore, we fuse the memory readout with multi-scale features from instance segmentation decoder through a carefully-designed multi-path fusion block. Extensive experiments conducted on both single-object and multi-object benchmarks demonstrate the effectiveness of the proposed method.

In addition to working towards more superior segmentation performance, another line of work [35, 62, 63] also explore the efficient memory storage to improve the inference efficiency of memory-based methods. Therefore, ISVOS can be combined with these approaches to develop both accurate and efficient VOS models.

## References

[1] Anurag Arnab and Philip HS Torr. Pixelwise instance segmentation with a dynamically instantiated network. In *CVPR*, 2017. 3

[2] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *ICCV*, 2019. 3

[3] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *CVPR*, 2017. 2

[4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 3

[5] Siddhartha Chandra and Iasonas Kokkinos. Fast, exact and multi-scale inference for semantic image segmentation with deep gaussian crfs. In *ECCV*, 2016. 8, 11

[6] 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. 5

[7] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blendmask: Top-down meets bottom-up for instance segmentation. In *CVPR*, 2020. 3

[8] Yuhua Chen, Jordi Pont-Tuset, Alberto Montes, and Luc Van Gool. Blazingly fast video object segmentation with pixel-wise metric learning. In *CVPR*, 2018. 2

[9] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 3

[10] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 2, 3, 4, 5, 7

[11] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*, 2021. 2, 3

[12] 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 *CVPR*, 2020. 5

[13] Ho Kei Cheng and Alexander G. Schwing. XMem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *ECCV*, 2022. 1, 2, 4, 5, 6, 7, 8, 9, 11, 12

[14] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion. In *CVPR*, 2021. 1, 2, 5, 6, 8

[15] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. In *NeurIPS*, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 9, 11, 12

[16] Jingchun Cheng, Yi-Hsuan Tsai, Shengjin Wang, and Ming-Hsuan Yang. Segflow: Joint learning for video object segmentation and optical flow. In *ICCV*, 2017. 2

[17] Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad El-badrawy, Ahsan Lodhi, and Harinandan Katam. Blenderproc. *arXiv preprint arXiv:1911.01911*, 2019. 5

[18] Brendan Duke, Abdalla Ahmed, Christian Wolf, Parham Aarabi, and Graham W Taylor. Sstvos: Sparse spatiotemporal transformers for video object segmentation. In *CVPR*, 2021. 2, 8

[19] Yuxin Fang, Shusheng Yang, Xinggang Wang, Yu Li, Chen Fang, Ying Shan, Bin Feng, and Wenyu Liu. Instances as queries. In *ICCV*, 2021. 2, 3

[20] Naiyu Gao, Yanhu Shan, Yupei Wang, Xin Zhao, Yinan Yu, Ming Yang, and Kaiqi Huang. Ssap: Single-shot instance segmentation with affinity pyramid. In *ICCV*, 2019. 3

[21] Melvyn A Goodale and A David Milner. Separate visual pathways for perception and action. *Trends in neurosciences*, 1992. 3

[22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 3

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 4

[24] 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 *CVPR*, 2021. 1, 2

[25] Yuan-Ting Hu, Jia-Bin Huang, and Alexander Schwing. Maskrnn: Instance level video object segmentation. In *NeurIPS*, 2017. 2

[26] Yuan-Ting Hu, Jia-Bin Huang, and Alexander G Schwing. Videomatch: Matching based video object segmentation. In *ECCV*, 2018. 2

[27] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring r-cnn. In *CVPR*, 2019. 3

[28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014. 5

[29] Shu Kong and Charless C Fowlkes. Recurrent pixel embedding for instance grouping. In *CVPR*, 2018. 3

[30] Mingxing Li, Li Hu, Zhiwei Xiong, Bang Zhang, Pan Pan, and Dong Liu. Recurrent dynamic embedding for video object segmentation. In *CVPR*, 2022. 5, 6, 7

[31] Xiaoxiao Li and Chen Change Loy. Video object segmentation with joint re-identification and attention-aware mask propagation. In *ECCV*, 2018. 2

[32] 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 *CVPR*, 2020. 5

[33] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017. 3

[34] Yu Li, Zhuoran Shen, and Ying Shan. Fast video object segmentation using the global context module. In *ECCV*, 2020. 1, 2

[35] Yongqing Liang, Xin Li, Navid Jafari, and Jim Chen. Video object segmentation with adaptive feature bank and uncertain-region refinement. In *NeurIPS*, 2020. 1, 2, 5, 8, 9, 12

[36] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 3, 5, 7

[37] Shu Liu, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *ICCV*, 2017. 3

[38] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018. 3

[39] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2017. 5

[40] Xiankai Lu, Wenguan Wang, Martin Danelljan, Tianfei Zhou, Jianbing Shen, and Luc Van Gool. Video object segmentation with episodic graph memory networks. In *ECCV*, 2020. 1, 2

[41] Yunyao Mao, Ning Wang, Wengang Zhou, and Houqiang Li. Joint inductive and transductive learning for video object segmentation. In *ICCV*, 2021. 8, 9, 12

[42] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *3DV*, 2016. 5

[43] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *CVPR*, 2019. 3

[44] Alejandro Newell, Zhiao Huang, and Jia Deng. Associative embedding: End-to-end learning for joint detection and grouping. In *NeurIPS*, 2017. 3

[45] Seoung Wug Oh, Joon-Young Lee, Kalyan Sunkavalli, and Seon Joo Kim. Fast video object segmentation by reference-guided mask propagation. In *CVPR*, 2018. 2

[46] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *ICCV*, 2019. 1, 2, 3, 4, 5, 7, 8, 9, 12

[47] Kwanyong Park, Sanghyun Woo, Seoung Wug Oh, In So Kweon, and Joon-Young Lee. Per-clip video object segmentation. In *CVPR*, 2022. 1

[48] Federico Perazzi, Anna Khoreva, Rodrigo Benenson, Bernt Schiele, and Alexander Sorkine-Hornung. Learning video object segmentation from static images. In *CVPR*, 2017. 2

[49] 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 *CVPR*, 2016. 2, 6

[50] 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. 2, 5, 6, 8

[51] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 3

[52] Hongje Seong, Junhyuk Hyun, and Euntai Kim. Kernelized memory network for video object segmentation. In *ECCV*, 2020. 2, 5

[53] Hongje Seong, Seoung Wug Oh, Joon-Young Lee, Seongwon Lee, Suhyeon Lee, and Euntai Kim. Hierarchical memory matching network for video object segmentation. In *ICCV*, 2021. 1, 2, 6

[54] Jianping Shi, Qiong Yan, Li Xu, and Jiaya Jia. Hierarchical image saliency detection on extended cssd. *TPAMI*, 2015. 5, 9, 12

[55] Yi-Hsuan Tsai, Ming-Hsuan Yang, and Michael J Black. Video segmentation via object flow. In *CVPR*, 2016. 2

[56] 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 *CVPR*, 2019. 2

[57] Haochen Wang, Xiaolong Jiang, Haibing Ren, Yao Hu, and Song Bai. Swiftnet: Real-time video object segmentation. In *CVPR*, 2021. 1, 2

[58] 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 *CVPR*, 2017. 5

[59] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *ECCV*, 2020. 3

[60] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. In *NeurIPS*, 2020. 3

[61] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018. 4

[62] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Shengping Zhang, and Wenxiu Sun. Efficient regional memory network for video object segmentation. In *CVPR*, 2021. 1, 2, 9, 12

[63] Kai Xu and Angela Yao. Accelerating video object segmentation with compressed video. In *CVPR*, 2022. 2, 9

[64] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. In *ECCV*, 2018. 2, 5, 6, 8

[65] Xiaohao Xu, Jinglu Wang, Xiao Li, and Yan Lu. Reliable propagation-correction modulation for video object segmentation. In *AAAI*, 2022. 5, 6

[66] Yu-Syuan Xu, Tsu-Jui Fu, Hsuan-Kung Yang, and Chun-Yi Lee. Dynamic video segmentation network. In *CVPR*, 2018. 2

[67] Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos K Katsaggelos. Efficient video object segmentation via network modulation. In *CVPR*, 2018. 2

[68] Shusheng Yang, Yuxin Fang, Xinggang Wang, Yu Li, Ying Shan, Bin Feng, and Wenyu Liu. Tracking instances as queries. *arXiv preprint arXiv:2106.11963*, 2021. 2

[69] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by foreground-background integration. In *ECCV*, 2020. 2, 8, 11

[70] Zongxin Yang, Yunchao Wei, and Yi Yang. Associating objects with transformers for video object segmentation. In *NeurIPS*, 2021. 5, 6, 9, 12

[71] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by multi-scale foreground-background integration. *TPAMI*, 2021. 8

[72] Zongxin Yang and Yi Yang. Decoupling features in hierarchical propagation for video object segmentation. In *NeurIPS*, 2022. 5, 6

[73] Yi Zeng, Pingping Zhang, Jianming Zhang, Zhe Lin, and Huchuan Lu. Towards high-resolution salient object detection. In *ICCV*, 2019. 5

[74] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable {detr}: Deformable transformers for end-to-end object detection. In *ICLR*, 2021. 4

## A. Multi-scale Inference.

Multi-scale evaluation is a commonly used trick in segmentation tasks [5, 13, 15] to boost the performance by merging the results of inputs under different data augmentations. Here we follow XMem [13] to apply image scaling and vertical mirroring and simply average the output probabilities to obtain the final masks.

| Method | MS | D16 val | | | D17 val | | |
|---|---|---|---|---|---|---|---|
| | | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
| CFBI [69] | ✓ | 90.7 | 89.6 | 91.7 | 83.3 | 80.5 | 86.0 |
| XMem [13] | ✓ | 92.7 | 92.0 | 93.5 | 88.2 | 85.4 | 91.0 |
| Ours | ✗ | 92.6 | 91.5 | 93.7 | 87.1 | 83.7 | 90.5 |
| Ours | ✓ | 92.9 | 92.2 | 93.6 | 88.6 | 85.8 | 91.4 |
| Ours* | ✗ | 92.8 | 91.8 | 93.8 | 88.2 | 84.5 | 91.9 |
| Ours* | ✓ | 93.4 | 92.5 | 94.2 | 89.8 | 86.7 | 93.0 |

Table 8. Results on DAVIS 2017 validation and YouTube-VOS validation split with different training data. D: DAVIS 2017, Y: YouTube 2019, S: static images, B: BL30K. ‡ denotes pretraining on the combined DAVIS and YouTube-VOS data.

The results in Table 8 imply that multi-scale inference improves the performance of ISVOS by 0.3% and 1.7% in terms of $\mathcal{J}\&\mathcal{F}$ on DAVIS 2016 / 2017 validation split, and ISVOS still outperforms existing methods.

## B. Results on Long video datasets

In order to further evaluate the long-term performance of ISVOS, we additionally test our method on the Long-time Video dataset [35], which contains three videos with more than 7,000 frames in total for validation. Considering the video duration is longer and the target object(s) will undergo distinct appearance deformation or scale variations, we set the maximum memory size to 64 during inference. The comparison results are shown in Table 9.

| Method | Long-time Video | | |
|---|---|---|---|
| | $\mathcal{J}\&\mathcal{F}$ | $\mathcal{J}$ | $\mathcal{F}$ |
| RMNet [62] | 59.8 | 59.7 | 60.0 |
| JOINT [41] | 67.1 | 64.5 | 69.6 |
| STM [46] | 80.6 | 79.9 | 81.3 |
| HMMN [54] | 81.5 | 79.9 | 83.0 |
| STCN [15] | 87.3 | 85.4 | 89.2 |
| AOT [70] | 84.3 | 83.2 | 85.4 |
| AFB-URR [35] | 83.7 | 82.9 | 84.5 |
| XMem [13] | 89.8 | 88.0 | 91.6 |
| Ours | **90.0** | **88.3** | **91.7** |

Table 9. Results on the Long-time Video dataset [35].

We can observe that ISVOS again achieves the best segmentation results measured in different metrics. It is worth mentioning that ISVOS beats the methods specifically designed for long videos, *e.g.*, AFB-URR [35] and XMem [13]. We believe the performance gain is resulted from taking advantage of the instance information in query frame to facilitate the semantic matching.

## C. More Visualizations

We show the predicted segmentation masks of ISVOS on DAVIS 2017 val, YouTube-VOS 2018 val, and Long-time Video dataset in Figure 9, Figure 10, Figure 11, respectively. For the short video datasets, *i.e.*, DAVIS and YouTube-VOS, the time interval is 5, while for the long video dataset, *i.e.*, Long-time Video dataset, the time interval is 1 since it is sparsely annotated. We can see that our method could generate accurate masks even for the objects with remarkable appearance variations.
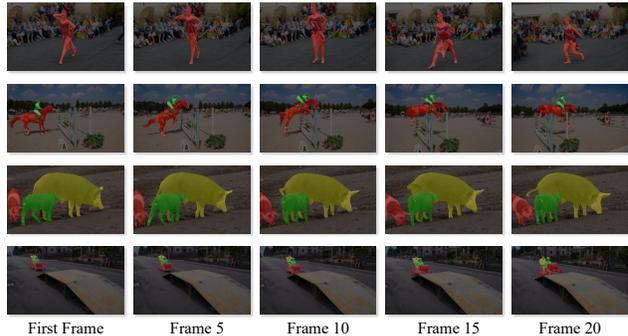


Figure 9. Segmentation results on DAVIS 2017 validation split.
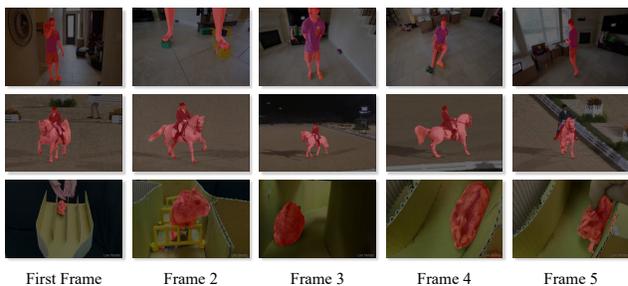


Figure 10. Segmentation results on YouTube-VOS 2018 validation split.



Figure 11. Segmentation results on Long-time Video dataset.