

Tracking Anything with Decoupled Video Segmentation

Ho Kei Cheng^{1†} Seoung Wug Oh² Brian Price² Alexander Schwing¹ Joon-Young Lee²
¹University of Illinois Urbana-Champaign ²Adobe Research
 {hokeikc2, aschwing}@illinois.edu, {seoh, bprice, jolee}@adobe.com



Figure 1. Visualization of our semi-online video segmentation results. Top: our algorithm (DEVA) extends Segment Anything (SAM) [27] to video for open-world video segmentation with no user input required. Bottom: DEVA performs text-prompted video segmentation for novel objects (with prompt “beyblade”, a type of spinning-top toy) by integrating Grounding-DINO [34] and SAM [27].

Abstract

Training data for video segmentation are expensive to annotate. This impedes extensions of end-to-end algorithms to new video segmentation tasks, especially in large-vocabulary settings. To ‘track anything’ without training on video data for every individual task, we develop a **decoupled video segmentation approach (DEVA)**, composed of task-specific image-level segmentation and class/task-agnostic bi-directional temporal propagation. Due to this design, we only need an image-level model for the target task (which is cheaper to train) and a universal temporal propagation model which is trained once and generalizes across tasks. To effectively combine these two modules, we use bi-directional propagation for (semi-)online fusion of segmentation hypotheses from different frames to generate a coherent segmenta-

tion. We show that this decoupled formulation compares favorably to end-to-end approaches in several data-scarce tasks including large-vocabulary video panoptic segmentation, open-world video segmentation, referring video segmentation, and unsupervised video object segmentation. Code is available at: [hkchengrex.github.io/Tracking-Anything-with-DEVA](https://github.com/hkchengrex/Tracking-Anything-with-DEVA).

1. Introduction

Video segmentation aims to segment and associate objects in a video. It is a fundamental task in computer vision and is crucial for many video understanding applications.

Most existing video segmentation approaches train end-to-end video-level networks on annotated video datasets. They have made significant strides on common benchmarks like YouTube-VIS [61] and Cityscape-VPS [24]. However,

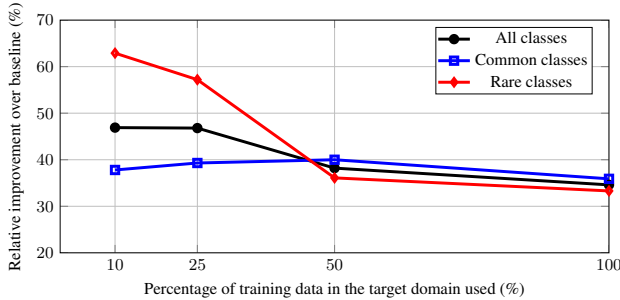


Figure 2. We plot relative $\overline{\text{VPQ}}$ increase of our decoupled approach over the end-to-end baseline when we vary the training data in the target domain (VIPSeg [39]). Common/rare classes are the top/bottom 50% most annotated object category in the training set. Our improvement is most significant ($>60\%$) in rare classes when there is a small amount of training data. This is because our decoupling allows the use of external class-agnostic temporal propagation data – data that cannot be used by existing end-to-end baselines. Details in Section 4.5.1.

these datasets have small vocabularies: YouTube-VIS contains 40 object categories, and Cityscape-VPS only has 19. It is questionable whether recent end-to-end paradigms are scalable to large-vocabulary, or even open-world video data. A recent larger vocabulary (124 classes) video segmentation dataset, VIPSeg [39], has been shown to be more difficult – using the same backbone, a recent method [30] achieves only 26.1 VPQ compared with 57.8 VPQ on Cityscape-VPS. To the best of our knowledge, recent video segmentation methods [2, 35] developed for the open-world setting (e.g., BURST [2]) are not end-to-end and are based on tracking of per-frame segmentation – further highlighting the difficulty of end-to-end training on large-vocabulary datasets. As the number of classes and scenarios in the dataset increases, it becomes more challenging to train and develop end-to-end video models to jointly solve segmentation and association, especially if annotations are scarce.

In this work, we aim to reduce reliance on the amount of target training data by leveraging external data *outside of the target domain*. For this, we propose to study *decoupled video segmentation*, which combines task-specific image-level segmentation and task-agnostic temporal propagation. Due to this design, we only need an image-level model for the target task (which is cheaper) and a universal temporal propagation model which is trained once and generalizes across tasks. Universal promptable image segmentation models like ‘segment anything’ (SAM) [27] have recently become available and serve as excellent candidates for the image-level model in a ‘track anything’ pipeline – Figure 1 shows some promising results of our integration with these methods.

Researchers have studied decoupled formulations before, as ‘tracking-by-detection’ [23, 51, 3]. However, these approaches often consider image-level detections im-

mutable, while the temporal model only associates detected objects. This formulation depends heavily on the quality of per-image detections and is sensitive to image-level errors.

In contrast, we develop a (semi-)online bi-directional propagation algorithm to 1) denoise image-level segmentation with in-clip consensus (Section 3.2.1), and 2) combine results from temporal propagation and in-clip consensus gracefully (Section 3.2.2). This bi-directional propagation allows temporally more coherent and potentially better results than those of an image-level model (see Figure 2).

We do not aim to replace end-to-end video approaches. Indeed, we emphasize that specialized frameworks on video tasks with sufficient video-level training data (e.g., YouTubeVIS [61]) outperform the developed method. Instead, we show that our decoupled approach acts as a strong baseline when an image model is available but video data is scarce. This is in spirit similar to pretraining of large language models [46]: a *task-agnostic* understanding of natural language is available before being finetuned on specific tasks – in our case, we learn propagation of segmentations of *class-agnostic* objects in videos via a temporal propagation module and make technical strides in applying this knowledge to specific tasks. The proposed decoupled approach transfers well to large-scale or open-world datasets, and achieves state-of-the-art results in large-scale video panoptic segmentation (VIPSeg [39]) and open-world video segmentation (BURST [2]). It also performs competitively on referring video segmentation (Ref-YouTubeVOS [49], Ref-DAVIS [22]) and unsupervised video object segmentation (DAVIS-16/17[5]) without end-to-end training.

To summarize:

- We propose using decoupled video segmentation that leverages external data, which allows it to generalize better to target tasks with limited annotations than end-to-end video approaches and allows us to seamlessly incorporate existing universal image segmentation models like SAM [27].
- We develop bi-directional propagation that denoises image segmentations and merges image segmentations with temporally propagated segmentations gracefully.
- We empirically show that our approach achieves favorable results in several important tasks including large-scale video panoptic segmentation, open-world video segmentation, referring video segmentation, and unsupervised video object segmentation.

2. Related Works

End-to-End Video Segmentation. Recent end-to-end video segmentation approaches [44, 21, 54, 4, 6, 14, 13] have made significant progress in tasks like Video Instance Segmentation (VIS) and Video Panoptic Segmentation (VPS), especially in closed and small vocabulary datasets like YouTube-VIS [61] and Cityscape-VPS [24].

However, these methods require end-to-end training and their scalability to larger vocabularies, where video data and annotations are expensive, is questionable. MaskProp [4] uses mask propagation to provide temporal information, but still needs to be trained end-to-end on the target task. This is because their mask propagation is not class-agnostic. We circumvent this training requirement and instead decouple the task into image segmentation and temporal propagation, each of which is easier to train with image-only data and readily available class-agnostic mask propagation data respectively.

Open-World Video Segmentation. Recently, an open-world video segmentation dataset BURST [2] has been proposed. It contains 482 object classes in diverse scenarios and evaluates open-world performance by computing metrics for the common classes (78, overlap with COCO [33]) and uncommon classes (404) separately. The baseline in BURST [2] predicts a set of object proposals using an image instance segmentation model trained on COCO [33] and associates the proposals frame-by-frame using either box IoU or STCN [11]. OWTB [35] additionally associates proposals using optical flow and pre-trained Re-ID features. Differently, we use bi-directional propagation that generates segmentations instead of simply associating existing segmentations – this reduces sensitivity to image segmentation errors. UVO [17] is another open-world video segmentation dataset and focuses on human actions. We mainly evaluate on BURST [2] as it is much more diverse and allows separate evaluation for common/uncommon classes.

Decoupled Video Segmentation. ‘Tracking-by-detection’ approaches [23, 51, 3] often consider image-level detections immutable and use a short-term temporal tracking model to associate detected objects. This formulation depends heavily on the quality of per-image detections and is sensitive to image-level errors. Related long-term temporal propagation works exist [19, 18], but they consider a single task and do not filter the image-level segmentation. We instead propose a general framework, with a bi-directional propagation mechanism that denoises the image segmentations and allows our result to potentially perform better than the image-level model.

Video Object Segmentation. Semi-supervised Video Object Segmentation (VOS) aims to propagate an initial ground-truth segmentation through a video [41, 40, 62, 9]. However, it does not account for any errors in the initial segmentation, and cannot incorporate new segmentation given by the image model at later frames. SAM-PT [47] combines point tracking with SAM [12] to create a video object segmentation pipeline, while our method tracks masks directly. We find a recent VOS algorithm [9] works well for our temporal propagation model. Our proposed bi-directional propagation is essential for bringing image segmentation models

and propagation models together as a unified video segmentation framework.

Unified Video Segmentation. Recent Video-K-Net [30] uses a unified framework for multiple video tasks but requires separate end-to-end training for each task. Unicorn [58], TarViS [1], and UNINEXT [59] share model parameters for different tasks, and train on all the target tasks end-to-end. They report lower tracking accuracy for objects that are not in the target tasks during training compared with class-agnostic VOS approaches, which might be caused by joint learning with class-specific features. In contrast, we only train an image segmentation model for the target task, while the temporal propagation model is always fully class-agnostic for generalization across tasks.

Segmenting/Tracking Anything. Concurrent to our work, Segment Anything (SAM) [27] demonstrates the effectiveness and generalizability of large-scale training for universal image segmentation, serving as an important foundation for open-world segmentation. Follow-up works [60, 12] extend SAM to video data by propagating the masks generated by SAM with video object segmentation algorithms. However, they rely on single-frame segmentation and lack the denoising capability of our proposed in-clip consensus approach.

3. Decoupled Video Segmentation

3.1. Formulation

Decoupled Video Segmentation. Our decoupled video segmentation approach is driven by an image segmentation model and a universal temporal propagation model. The image model, trained specifically on the target task, provides task-specific image-level segmentation hypotheses. The temporal propagation model, trained on class-agnostic mask propagation datasets, associates and propagates these hypotheses to segment the whole video. This design separates the learning of task-specific segmentation and the learning of general video object segmentation, leading to a robust framework even when data in the target domain is scarce and insufficient for end-to-end learning.

Notation. Using t as the time index, we refer to the corresponding frame and its final segmentation as I_t and \mathbf{M}_t respectively. In this paper, we represent a segmentation as a set of non-overlapping per-object binary segments, *i.e.*, $\mathbf{M}_t = \{m_i, 0 < i \leq |\mathbf{M}_t|\}$, where $m_i \cap m_j = \emptyset$ if $i \neq j$.

The image segmentation model $\text{Seg}(I)$ takes an image I as input and outputs a segmentation. We denote its output segmentation at time t as $\text{Seg}(I_t) = \text{Seg}_t = \{s_i, 0 < i \leq |\text{Seg}_t|\}$, which is also a set of non-overlapping binary segments. This segmentation model can be swapped for different target tasks, and users can be in the loop to correct the segmentation as we do not limit its internal architecture.

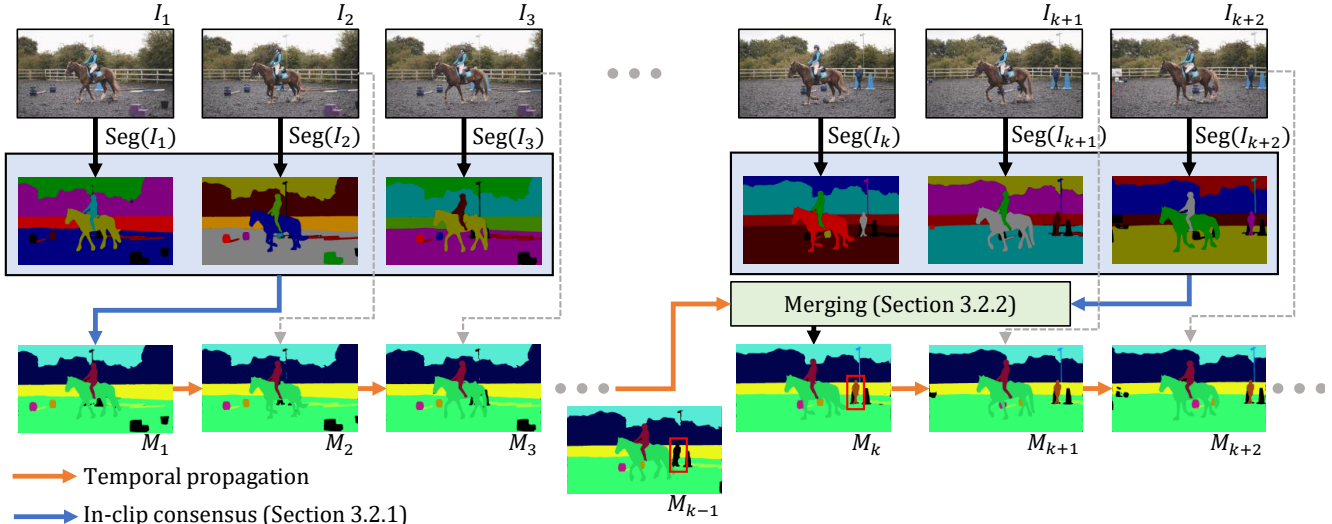


Figure 3. Overview of our framework. We first filter image-level segmentations with in-clip consensus (Section 3.2.1) and temporally propagate this result forward. To incorporate a new image segmentation at a later time step (for previously unseen objects, e.g., red box), we merge the propagated results with in-clip consensus as described in Section 3.2.2. Specifics of temporal propagation are in the appendix.

The temporal propagation model $\text{Prop}(\mathbf{H}, I)$ takes a collection of segmented frames (memory) \mathbf{H} and a query image I as input and segments the query frame with the objects in the memory. For instance, $\text{Prop}(\{I_1, \mathbf{M}_1\}, I_2)$ propagates the segmentation \mathbf{M}_1 from the first frame I_1 to the second frame I_2 . Unless mentioned explicitly, the memory \mathbf{H} contains all past segmented frames.

Overview. Figure 3 illustrates the overall pipeline. At a high level, we aim to propagate segmentations discovered by the image segmentation model to the full video with temporal propagation. We mainly focus on the (semi-)online setting. Starting from the first frame, we use the image segmentation model for initialization. To denoise errors from single-frame segmentation, we look at a small clip of a few frames in the near future (in the online setting, we only look at the current frame) and reach an in-clip consensus (Section 3.2.1) as the output segmentation. Afterward, we use the temporal propagation model to propagate the segmentation to subsequent frames. We modify an off-the-shelf state-of-the-art video object segmentation XMem [9] as our temporal propagation model, with details given in the appendix. The propagation model itself cannot segment new objects that appear in the scene. Therefore, we periodically incorporate new image segmentation results using the same in-clip consensus as before and merge the consensus with the propagated result (Section 3.2.2). This pipeline combines the strong temporal consistency from the propagation model (past) and the new semantics from the image segmentation model (future), hence the name *bi-directional propagation*. Next, we will discuss the bi-directional propagation pipeline in detail.

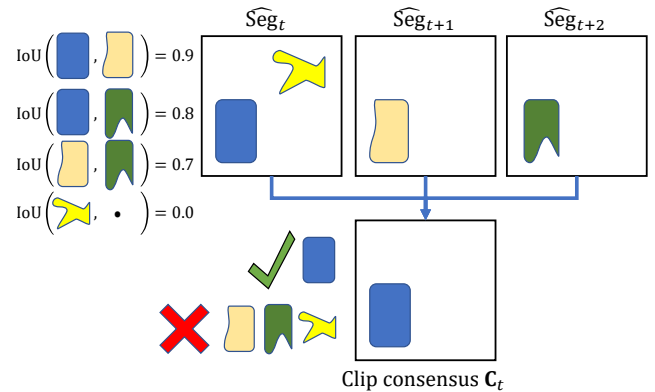


Figure 4. A simple illustration of in-clip consensus. The top three squares represent object proposals from three different frames aligned to time t . The blue shape is the most supported by other object proposals and is selected as output. The yellow shape is not supported by any and is ruled out as noise. The remaining are not used due to significant overlap with the selected (blue) shape.

3.2. Bi-Directional Propagation

3.2.1 In-clip Consensus

Formulation. In-clip consensus operates on the image segmentations of a small future clip of n frames ($\text{Seg}_t, \text{Seg}_{t+1}, \dots, \text{Seg}_{t+n-1}$) and outputs a denoised consensus \mathbf{C}_t for the current frame. In the online setting, $n = 1$ and $\mathbf{C}_t = \text{Seg}_t$. In the subsequent discussion, we focus on the semi-online setting, as consensus computation in the online setting is straightforward. As an overview, we first obtain a set of *object proposals* on the target frame t via spatial alignment, merge the object proposals into a combined rep-

resentation in a second step, and optimize for an indicator variable to choose a subset of proposals as the output in an integer program. Figure 4 illustrates this in-clip consensus computation in a stylized way and we provide details regarding each of the three aforementioned steps (spatial alignment, representation, and integer programming) next.

Spatial Alignment. As the segmentations $(\text{Seg}_t, \text{Seg}_{t+1}, \dots, \text{Seg}_{t+n-1})$ correspond to different time steps, they might be spatially misaligned. This misalignment complicates the computation of correspondences between segments. To align segmentations Seg_{t+i} with frame t , techniques like optical flow warping are applicable. In this paper, we simply re-use the temporal propagation model to find the aligned segmentation $\widehat{\text{Seg}}_{t+i}$ (note $\widehat{\text{Seg}}_t = \text{Seg}_t$) via

$$\widehat{\text{Seg}}_{t+i} = \text{Prop}(\{I_{t+i}, \text{Seg}_{t+i}\}, I_t), 0 < i < n. \quad (1)$$

Note, the propagation model here only uses one frame as memory at a time and this temporary memory $\{I_{t+i}, \text{Seg}_{t+i}\}$ is discarded immediately after alignment. It does not interact with the global memory \mathbf{H} .

Representation. Recall that we represent a segmentation as a set of non-overlapping per-object binary segments. After aligning all the segmentations to frame t , each segment is an *object proposal* for frame I_t . We refer to the union of all these proposals via \mathbf{P} (time index omitted for clarity):

$$\mathbf{P} = \bigcup_{i=0}^{n-1} \widehat{\text{Seg}}_{t+i} = \{p_i, 0 < i \leq |\mathbf{P}|\}. \quad (2)$$

The output of consensus voting is represented by an indicator variable $v^* \in \{0, 1\}^{|\mathbf{P}|}$ that combines segments into the consensus output \mathbf{C}_t :

$$\mathbf{C}_t = \{p_i | v_i^* = 1\} = \{c_i, 0 < i \leq |\mathbf{C}|\}. \quad (3)$$

We resolve overlapping segments c_i in \mathbf{C}_t by prioritizing smaller segments as they are more vulnerable to being majorly displaced by overlaps. This priority is implemented by sequentially rendering the segments c_i on an image in descending order of area. We optimize for v based on two simple criteria:

1. Lone proposals p_i are likely to be noise and should not be selected. Selected proposals should be supported by other (unselected) proposals.
2. Selected proposals should not overlap significantly with each other.

We combine these criteria in an integer programming problem which we describe next.

Integer Programming. We aim to optimize the indicator variable v to achieve the above two objectives, by addressing the following integer programming problem:

$$v^* = \underset{v}{\text{argmax}} \sum_i (\text{Supp}_i + \text{Penal}_i) \text{ s.t. } \sum_{i,j} \text{Overlap}_{ij} = 0. \quad (4)$$

Next, we discuss each of the terms in the program in detail.

First, we define the pairwise Intersection-over-Union (IoU) between the i -th proposal and the j -th proposal as:

$$\text{IoU}_{ij} = \text{IoU}_{ji} = \frac{|p_i \cap p_j|}{|p_i \cup p_j|}, 0 \leq \text{IoU}_{ij} \leq 1. \quad (5)$$

The i -th proposal *supports* the j -th proposal if $\text{IoU}_{ij} > 0.5$ – the higher the IoU, the stronger the support. The more support a segment has, the more favorable it is to be selected. To maximize the total support of selected segments, we maximize the below objective for all i :

$$\text{Supp}_i = v_i \sum_j \begin{cases} \text{IoU}_{ij}, & \text{if } \text{IoU}_{ij} > 0.5 \text{ and } i \neq j \\ 0, & \text{otherwise} \end{cases}. \quad (6)$$

Additionally, proposals that support each other should not be selected together as they significantly overlap. This is achieved by constraining the following term to zero:

$$\text{Overlap}_{ij} = \begin{cases} v_i v_j, & \text{if } \text{IoU}_{ij} > 0.5 \text{ and } i \neq j \\ 0, & \text{otherwise} \end{cases}. \quad (7)$$

Lastly, we introduce a penalty for selecting any segment for 1) tie-breaking when a segment has no support, and 2) excluding noisy segments, with weight α :

$$\text{Penal}_i = -\alpha v_i. \quad (8)$$

We set the tie-breaking weight $\alpha = 0.5$. For all but the first frame, we merge \mathbf{C}_t with the propagated segmentation $\text{Prop}(\mathbf{H}, I_t)$ into the final output \mathbf{M}_t as described next.

3.2.2 Merging Propagation and Consensus

Formulation. Here, we seek to merge the propagated segmentation $\text{Prop}(\mathbf{H}, I_t) = \mathbf{R}_t = \{r_i, 0 < i \leq |\mathbf{R}|\}$ (from the past) with the consensus $\mathbf{C}_t = \{c_j, 0 < j \leq |\mathbf{C}|\}$ (from the near future) into a single segmentation \mathbf{M}_t . We associate segments from these two segmentations and denote the association with an indicator a_{ij} which is 1 if r_i associates with c_j , and 0 otherwise. Different from the in-clip consensus, these two segmentations contain fundamentally different information. Thus, we do not eliminate any segments and instead fuse all pairs of associated segments while letting the unassociated segments pass through to the output. Formally, we obtain the final segmentation via

$$\mathbf{M}_t = \{r_i \cup c_j | a_{ij} = 1\} \cup \{r_i | \forall_j a_{ij} = 0\} \cup \{c_j | \forall_i a_{ij} = 0\}, \quad (9)$$

where overlapping segments are resolved by prioritizing the smaller segments as discussed in Section 3.2.1.

Maximizing Association IoU. We find a_{ij} by maximizing the pairwise IoU of all associated pairs, with a minimum association IoU of 0.5. This is equivalent to a maximum bipartite matching problem, with r_i and c_j as vertices and edge weight e_{ij} given by

$$e_{ij} = \begin{cases} \text{IoU}(r_i, c_j), & \text{if } \text{IoU}(r_i, c_j) > 0.5 \\ -1, & \text{otherwise} \end{cases}. \quad (10)$$

Requiring any matched pairs from two non-overlapping segmentations to have $\text{IoU} > 0.5$ leads to a unique matching, as shown in [26]. Therefore, a greedy solution of setting $a_{ij} = 1$ if $e_{ij} > 0$ and 0 otherwise suffices to obtain an optimal result.

Segment Deletion. As an implementation detail, we delete inactive segments from the memory to reduce computational costs. We consider a segment r_i inactive when it fails to associate with any segments c_j from the consensus for consecutive L times. Such objects might have gone out of view or were a misdetection. Concretely, we associate a counter cnt_i with each propagated segment r_i , initialized as 0. When r_i is not associated with any segments c_j from the consensus, i.e., $\forall_j a_{ij} = 0$, we increment cnt_i by 1 and reset cnt_i to 0 otherwise. When cnt_i reaches the pre-defined threshold L , the segment r_i is deleted from the memory. We set $L = 5$ in all our experiments.

4. Experiments

We first present our main results using a large-scale video panoptic segmentation dataset (VIPSeg [39]) and an open-world video segmentation dataset (BRUST [2]). Next, we show that our method also works well for referring video object segmentation and unsupervised video object segmentation. We present additional results on the smaller-scale YouTubeVIS dataset in the appendix, but unsurprisingly recent end-to-end specialized approaches perform better because a sufficient amount of data is available in this case. Figure 1 visualizes some results of the integration of our approach with universal image segmentation models like SAM [27] or Grounding-Segment-Anything [34, 27]. By default, we merge in-clip consensus with temporal propagation every 5 frames with a clip size of $n = 3$ in the semi-online setting, and $n = 1$ in the online setting. We evaluate all our results using either official evaluation codebases or official servers. We use image models trained with standard training data for each task (using open-sourced models whenever available) and a universal temporal propagation module for all tasks unless otherwise specified.

The temporal propagation model is based on XMem [9], and is trained in a class-agnostic fashion with image segmentation datasets [50, 53, 63, 29, 8] and video object segmentation datasets [57, 41, 42]. With the long-term memory of XMem [9], our model can handle long videos with ease.

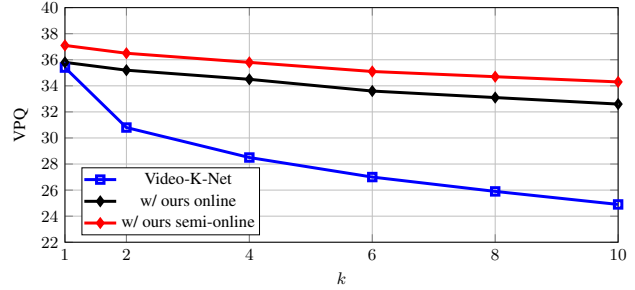


Figure 5. Performance trend comparison of Video-K-Net [30] and our decoupled approach with the same base model. Ours decreases slower with larger k , indicating that the proposed decoupled method has a better long-term propagation.

We use top-k filtering [10] with $k = 30$ following [9]. The performance of our modified propagation model on common video object segmentation benchmarks (DAVIS [41], YouTubeVOS [57], and MOSE [15]) are listed in the appendix.

4.1. Large-Scale Video Panoptic Segmentation

We are interested in addressing the large vocabulary setting. To our best knowledge, VIPSeg [39] is currently the largest scale in-the-wild panoptic segmentation dataset, with 58 things classes and 66 stuff classes in 3,536 videos of 232 different scenes.

Metrics. To evaluate the quality of the result, we adopt the commonly used VPQ (Video Panoptic Quality) [24] and STQ (Segmentation and Tracking Quality) [55] metrics. VPQ extends image-based PQ (Panoptic Quality) [26] to video data by matching objects in sliding windows of k frames (denoted VPQ^k). When $k = 1$, $\text{VPQ} = \text{PQ}$ and associations of segments between frames are ignored. Correct long-range associations, which are crucial for object tracking and video editing tasks, are only evaluated with a large value of k . For a more complete evaluation of VPS, we evaluate $k \in \{1, 2, 4, 6, 8, 10, \infty\}$. Note, VPQ^∞ considers the entire video as a tube and requires global association. We additionally report $\bar{\text{VPQ}}$, which is the average of VPQ^∞ and the arithmetic mean of $\text{VPQ}^{\{1,2,4,6,8,10\}}$. This weights VPQ^∞ higher as it represents video-level performance, while the other metrics only assess frame-level or clip-level results. STQ is proposed in STEP [55] and is the geometric mean of AQ (Association Quality) and SQ (Segmentation Quality). It evaluates pixel-level associations and semantic segmentation quality respectively. We refer readers to [24] and [55] for more details on VPQ and STQ.

Main Results. Table 1 summarizes our findings. To assess generality, we study three models as image segmentation input (PanoFCN [31], Mask2Former [7], and Video-K-Net [30]) to our decoupled approach. The weights of these image models are initialized by pre-training on the COCO panoptic dataset [33] and subsequently fine-tuned

Backbone				VPQ ¹	VPQ ²	VPQ ⁴	VPQ ⁶	VPQ ⁸	VPQ ¹⁰	VPQ [∞]	\overline{VPQ}	STQ
Clip-PanoFCN		end-to-end [39]	semi-online	27.3	26.0	24.2	22.9	22.1	21.5	18.1	21.1	28.3
Clip-PanoFCN		decoupled (ours)	online	29.5	28.9	28.1	27.2	26.7	26.1	25.0	26.4	35.7
Clip-PanoFCN		decoupled (ours)	semi-online	31.3	30.8	30.1	29.4	28.8	28.3	27.1	28.4	35.8
Video-K-Net	R50	end-to-end [30]	online	35.4	30.8	28.5	27.0	25.9	24.9	21.7	25.2	33.7
Video-K-Net	R50	decoupled (ours)	online	35.8	35.2	34.5	33.6	33.1	32.6	30.5	32.3	38.4
Video-K-Net	R50	decoupled (ours)	semi-online	37.1	36.5	35.8	35.1	34.7	34.3	32.3	33.9	38.6
Mask2Former	R50	decoupled (ours)	online	41.0	40.2	39.3	38.4	37.9	37.3	33.8	36.4	41.1
Mask2Former	R50	decoupled (ours)	semi-online	42.1	41.5	40.8	40.1	39.7	39.3	36.1	38.3	41.5
Video-K-Net	Swin-B	end-to-end [30]	online	49.8	45.2	42.4	40.5	39.1	37.9	32.6	37.5	45.2
Video-K-Net	Swin-B	decoupled (ours)	online	48.2	47.4	46.5	45.6	45.1	44.5	42.0	44.1	48.6
Video-K-Net	Swin-B	decoupled (ours)	semi-online	50.0	49.3	48.5	47.7	47.3	46.8	44.5	46.4	48.9
Mask2Former	Swin-B	decoupled (ours)	online	55.3	54.6	53.8	52.8	52.3	51.9	49.0	51.2	52.4
Mask2Former	Swin-B	decoupled (ours)	semi-online	56.0	55.4	54.6	53.9	53.5	53.1	50.0	52.2	52.2

Table 1. Comparisons of end-to-end approaches (e.g., state-of-the-art Video-K-Net [30]) with our decoupled approach on the large-scale video panoptic segmentation dataset VIPSeg [39]. Our method scales with better image models and performs especially well with large k where long-term associations are considered. All baselines are reproduced using official codebases.

Method		Validation			Test		
		OWTA _{all}	OWTA _{com}	OWTA _{unc}	OWTA _{all}	OWTA _{com}	OWTA _{unc}
Mask2Former	w/ Box tracker [2]	60.9	66.9	24.0	55.9	61.0	24.6
Mask2Former	w/ STCN tracker [2]	64.6	71.0	25.0	57.5	62.9	23.9
OWTB [35]		55.8	59.8	38.8	56.0	59.9	38.3
Mask2Former	w/ ours online	69.5	74.6	42.3	70.1	75.0	44.1
Mask2Former	w/ ours semi-online	69.9	75.2	41.5	70.5	75.4	44.1
EntitySeg	w/ ours online	68.8	72.7	49.6	69.5	72.9	53.0
EntitySeg	w/ ours semi-online	69.5	73.3	50.5	69.8	73.1	53.3

Table 2. Comparison to baselines in the open-world video segmentation dataset BURST [2]. ‘com’ stands for ‘common classes’ and ‘unc’ stands for ‘uncommon classes’. Our method performs better in both – in the common classes with Mask2Former [7] image backbone, and in the uncommon classes with EntitySeg [43]. The ability to switch image backbones is one of the main advantages of our decoupled formulation. Baseline performances are transcribed from [2].

on VIPSeg [39]. Our method outperforms both baseline Clip-PanoFCN [39] and state-of-the-art Video-K-Net [30] with the same backbone, especially if k is large, *i.e.*, when long-term associations are more important. Figure 5 shows the performance trend with respect to k . The gains for large values of k highlight the use of a decoupled formulation over end-to-end training: the latter struggles with associations eventually, as training sequences aren’t arbitrarily long. Without any changes to our generalized mask propagation module, using a better image backbone (*e.g.*, SwinB [36]) leads to noticeable improvements. Our method can likely be coupled with future advanced methods in image segmentation for even better performance.

4.2. Open-World Video Segmentation

Open-world video segmentation addresses the difficult problem of discovering, segmenting, and tracking objects

in the wild. BURST [2] is a recently proposed dataset that evaluates open-world video segmentation. It contains diverse scenarios and 2,414 videos in its validation/test sets. There are a total of 482 object categories, 78 of which are ‘common’ classes while the rest are ‘uncommon’.

Metrics. Following [2], we assess Open World Tracking Accuracy (OWTA), computed separately for ‘all’, ‘common’, and ‘uncommon’ classes. False positive tracks are not directly penalized in the metrics as the ground-truth annotations are not exhaustive for all objects in the scene, but indirectly penalized by requiring the output mask to be mutually exclusive. We refer readers to [2, 37] for details.

Main Results. Table 2 summarizes our findings. We study two image segmentation models: Mask2Former [7], and EntitySeg [43], both of which are pretrained on the COCO [33] dataset. The Mask2Former weight is trained for the instance segmentation task, while EntitySeg is trained

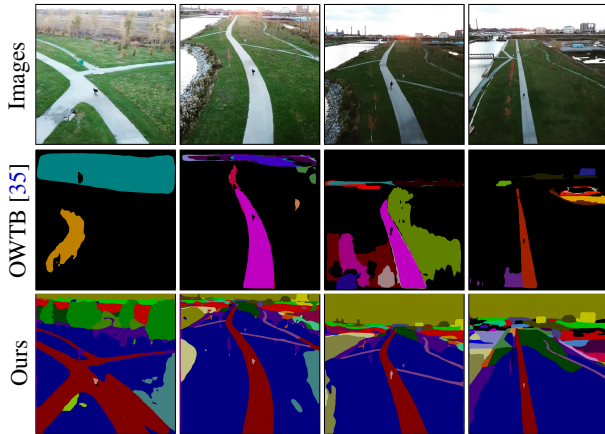


Figure 6. An in-the-wild result in the BURST [2] dataset. Note, we can even track the small skateboarder (pink mask on the road).

for ‘entity segmentation’, that is to segment all visual entities without predicting class labels. We find EntitySeg works better for novel objects, as it is specifically trained to do so. Being able to plug and play the latest development of open-world image segmentation models without any fine-tuning is one of the major advantages of our formulation.

Our approach outperforms the baselines, which all follow the ‘tracking-by-detection’ paradigm. In these baselines, segmentations are detected every frame, and a short-term temporal module is used to associate these segmentations between frames. This paradigm is sensitive to misdetections in the image segmentation model. ‘Box tracker’ uses per-frame object IoU; ‘STCN tracker’ uses a pretrained STCN [11] mask propagation network; and OWTB [35] uses a combination of IoU, optical flow, and Re-ID features. We also make use of mask propagation, but we go beyond the setting of simply associating existing segmentations – our bi-directional propagation allows us to improve upon the image segmentations and enable long-term tracking. Figure 6 compares our results on one of the videos in BURST to OWTB [35].

4.3. Referring Video Segmentation

Referring video segmentation takes a text description of an object as input and segments the target object. We experiment on Ref-DAVIS17 [22] and Ref-YouTubeVOS [49] which augments existing video object segmentation datasets [41, 57] with language expressions. Following [56], we assess \mathcal{J} & \mathcal{F} which is the average of Jaccard index (\mathcal{J}), and boundary F1-score (\mathcal{F}).

Table 3 tabulates our results. We use an image-level ReferFormer [56] as the image segmentation model. We find that the quality of referring segmentation has a high variance across the video (e.g., the target object might be too small at the beginning of the video). As in all competing approaches [49, 56, 16], we opt for an offline setting to

reduce this variance. Concretely, we perform the initial in-clip consensus by selecting 10 uniformly spaced frames in the video and using the frame with the highest confidence given by the image model as a ‘key frame’ for aligning the other frames. We then forward- and backward-propagate from the key frame without incorporating additional image segmentations. We give more details in the appendix. Our method outperforms other approaches.

Method	Ref-DAVIS [22]	Ref-YTVOS [49]
URVOS [49]	51.6	47.2
ReferFormer [56]	60.5	62.4
VLT [16]	61.6	63.8
Ours	66.3	66.0

Table 3. \mathcal{J} & \mathcal{F} comparisons on two referring video segmentation datasets. Ref-YTVOS stands for Ref-YouTubeVOS [49].

4.4. Unsupervised Video Object Segmentation

Unsupervised video object segmentation aims to find and segment salient target object(s) in a video. We evaluate on DAVIS-16 [41] (single-object) and DAVIS-17 [5] (multi-object). In the single-object setting, we use the image saliency model DIS [45] as the image model and employ an offline setting as in Section 4.3. In the multi-object setting, since the image saliency model only segments one object, we instead use EntitySeg [43] and follow our semi-online protocol on open-world video segmentation in Section 4.2. Table 4 summarizes our findings. Please refer to the appendix for details.

Method	D16-val	D17-val	D17-td
RTNet [48]	85.2	-	-
PMN [28]	85.9	-	-
UnOVOST [38]	-	67.9	58.0
Propose-Reduce [32]	-	70.4	-
Ours	88.9	73.4	62.1

Table 4. \mathcal{J} & \mathcal{F} comparisons on three unsupervised video object segmentation datasets: DAVIS16 validation (D16-val), DAVIS17 validation (D17-val), and DAVIS17 test-dev (D17-td). Missing entries mean that the method did not report results on that dataset.

4.5. Ablation Studies

4.5.1 Varying Training Data

Here, we vary the amount of training data in the target domain (VIPSeg [39]) to measure the sensitivity of end-to-end approaches vs. our decoupled approach. We subsample different percentages of videos from the training set

<i>Varying clip size</i>	VPQ ¹	VPQ ¹⁰	$\overline{\text{VPQ}}$	STQ	FPS
$n = 1$	41.0	37.3	36.4	41.1	10.3
$n = 2$	40.4	37.2	36.3	39.0	9.8
$n = 3$	42.1	39.3	38.3	41.5	7.8
$n = 4$	42.1	39.1	38.5	42.3	6.6
$n = 5$	41.7	38.9	38.3	42.8	5.6
<i>Varying merge freq.</i>	VPQ ¹	VPQ ¹⁰	$\overline{\text{VPQ}}$	STQ	FPS
Every 3 frames	42.2	39.2	38.4	42.6	5.2
Every 5 frames	42.1	39.3	38.3	41.5	7.8
Every 7 frames	41.5	39.0	35.7	40.5	8.4
<i>Spatial Align?</i>	VPQ ¹	VPQ ¹⁰	$\overline{\text{VPQ}}$	STQ	FPS
Yes	42.1	39.3	38.3	41.5	7.8
No	36.7	33.9	32.8	33.7	9.2

Table 5. Performances of our method on VIPSeg [39] with different hyperparameters and design choices. By default, we use a clip size of $n = 3$ and a merge frequency of every 5 frames with spatial alignment for a balance between performance and speed.

to train Video-K-Net-R50 [30] (all networks are still pre-trained with COCO-panoptic [33]). We then compare end-to-end performances with our (semi-online) decoupled performances (the temporal propagation model is unchanged as it does not use any data from the target domain). Figure 1 plots our findings – our model has a much higher relative $\overline{\text{VPQ}}$ improvement over the baseline Video-K-Net for rare classes if little training data is available.

4.5.2 In-Clip Consensus

Here we explore hyperparameters and design choices in in-clip consensus. Table 5 tabulates our performances with different *clip sizes*, different *frequencies* of merging in-clip consensus with temporal propagation, and whether to use *spatial alignment* during in-clip consensus. Mask2Former-R50 is used as the backbone in all entries. For clip size $n = 2$, tie-breaking is ambiguous. A large clip is more computationally demanding and potentially leads to inaccurate spatial alignment as the appearance gap between frames in the clip increases. A high merging frequency reduces the delay between the appearance of a new object and its detection in our framework but requires more computation. By default, we use a clip size $n = 3$, merge consensus with temporal propagation every 5 frames, and enable spatial alignment for a balance between performance and speed.

4.5.3 Using Temporal Propagation

Here, we compare different approaches for using temporal propagation in a decoupled setting. Tracking-by-detection approaches [23, 51, 3] typically detect segmentation at every frame and use temporal propagation to associate these per-frame segmentations. We test these short-term asso-

ciation approaches using 1) mask IoU between adjacent frames, 2) mask IoU of adjacent frames warped by optical flow from RAFT [52], and 3) query association [20] of query-based segmentation [7] between adjacent frames. We additionally compare with variants of our temporal propagation method: 4) ‘ShortTrack’, where we consider only short-term tracking by re-initializing the memory \mathbf{H} every frame, and 5) ‘TrustImageSeg’, where we explicitly trust the consensus given by the image segmentations over temporal propagation by discarding segments that are not associated with a segment in the consensus (i.e., dropping the middle term in Eq. (9)). Table 6 tabulates our findings. For all entries, we use Mask2Former-R50 [7] in the online setting on VIPSeg [39] for fair comparisons.

Temporal scheme	VPQ ¹	VPQ ⁴	VPQ ¹⁰	$\overline{\text{VPQ}}$	STQ
Mask IoU	39.9	32.7	27.7	27.6	34.5
Mask IoU+flow	40.2	33.7	28.8	28.6	37.0
Query assoc.	40.4	33.1	28.1	28.0	35.8
‘ShortTrack’	40.6	33.3	28.3	28.2	37.2
‘TrustImageSeg’	40.3	37.5	33.7	33.2	37.9
Ours, bi-direction	41.0	39.3	37.3	36.4	41.1

Table 6. Performances of different temporal schema on VIPSeg [39]. Our bi-directional propagation scheme is necessary for the final high performance.

4.6. Limitations

As the temporal propagation model is task-agnostic, it cannot detect new objects by itself. As shown by the red boxes in Figure 3, the new object in the scene is missing from \mathbf{M}_{k-1} and can only be detected in \mathbf{M}_k – this results in delayed detections relating to the frequency of merging with in-clip consensus. Secondly, we note that end-to-end approaches still work better when training data is sufficient, i.e., in smaller vocabulary settings like YouTubeVIS [61] as shown in the appendix. But we think decoupled methods are more promising in large-vocabulary/open-world settings.

5. Conclusion

We present **DEVA**, a decoupled video segmentation approach for ‘tracking anything’. It uses a bi-directional propagation technique that effectively scales image segmentation methods to video data. Our approach critically leverages external task-agnostic data to reduce reliance on the target task, thus generalizing better to tasks with scarce data than end-to-end approaches. Combined with universal image segmentation models, our decoupled paradigm demonstrates state-of-the-art performance as a first step towards open-world large-vocabulary video segmentation.

Acknowledgments. Work supported in part by NSF grants 2008387, 2045586, 2106825, MRI 1725729 (HAL [25]), and NIFA award 2020-67021-32799.

References

- [1] Ali Athar, Alexander Hermans, Jonathon Luiten, Deva Ramanan, and Bastian Leibe. Tarvis: A unified approach for target-based video segmentation. *arXiv preprint arXiv:2301.02657*, 2023.
- [2] Ali Athar, Jonathon Luiten, Paul Voigtlaender, Tarasha Khurana, Achal Dave, Bastian Leibe, and Deva Ramanan. Burst: A benchmark for unifying object recognition, segmentation and tracking in video. In *WACV*, 2023.
- [3] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, 2019.
- [4] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances in video with mask propagation. In *CVPR*, 2020.
- [5] Sergi Caelles, Jordi Pont-Tuset, Federico Perazzi, Alberto Montes, Kevis-Kokitsi Maninis, and Luc Van Gool. The 2019 davis challenge on vos: Unsupervised multi-object segmentation. In *arXiv preprint arXiv:1905.00737*, 2019.
- [6] Bowen Cheng, Anwesa Choudhuri and Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexander G. Schwing. Mask2former for video instance segmentation. In <https://arxiv.org/abs/2112.10764>, 2021.
- [7] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022.
- [8] 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.
- [9] Ho Kei Cheng and Alexander G Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In *ECCV*, 2022.
- [10] 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.
- [11] 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.
- [12] Yangming Cheng, Liulei Li, Yuanyou Xu, Xiaodi Li, Zongxin Yang, Wenguan Wang, and Yi Yang. Segment and track anything. In *arXiv preprint arXiv:2305.06558*, 2023.
- [13] A. Choudhuri, G. Chowdhary, and A. G. Schwing. Assignment-Space-Based Multi-Object Tracking and Segmentation. In *ICCV*, 2021.
- [14] A. Choudhuri, G. Chowdhary, and A. G. Schwing. Context-Aware Relative Object Queries to Unify Video Instance and Panoptic Segmentation. In *CVPR*, 2023.
- [15] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, Philip HS Torr, and Song Bai. MOSE: A new dataset for video object segmentation in complex scenes. In *ICCV*, 2023.
- [16] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vlt: Vision-language transformer and query generation for referring segmentation. In *TPAMI*, 2022.
- [17] Yuming Du, Wen Guo, Yang Xiao, and Vincent Lepetit. Uvo challenge on video-based open-world segmentation 2021: 1st place solution. *ICCV Workshop*, 2021.
- [18] Shubhika Garg and Vidit Goel. Mask selection and propagation for unsupervised video object segmentation. In *WACV*, 2021.
- [19] Vidit Goel, Jiachen Li, Shubhika Garg, Harsh Maheshwari, and Humphrey Shi. Msn: efficient online mask selection network for video instance segmentation. In *CVPR Workshop*, 2021.
- [20] De-An Huang, Zhiding Yu, and Anima Anandkumar. Minvis: A minimal video instance segmentation framework without video-based training. In *NeurIPS*, 2022.
- [21] Sukjun Hwang, Miran Heo, Seoung Wug Oh, and Seon Joo Kim. Video instance segmentation using inter-frame communication transformers. *NeurIPS*, 2021.
- [22] Anna Khoreva, Anna Rohrbach, and Bernt Schiele. Video object segmentation with language referring expressions. In *ACCV*, 2019.
- [23] Chanh Kim, Fuxin Li, Arridhana Ciptadi, and James M Rehg. Multiple hypothesis tracking revisited. In *ICCV*, 2015.
- [24] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Video panoptic segmentation. In *CVPR*, 2020.
- [25] Volodymyr Kindratenko, Dawei Mu, Yan Zhan, John Maloney, Sayed Hadi Hashemi, Benjamin Rabe, Ke Xu, Roy Campbell, Jian Peng, and William Gropp. Hal: Computer system for scalable deep learning. In *PEARC*, 2020.
- [26] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019.
- [27] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *arXiv preprint arXiv:2304.02643*, 2023.
- [28] Minhyeok Lee, Suhwan Cho, Seunghoon Lee, Chaewon Park, and Sangyoun Lee. Unsupervised video object segmentation via prototype memory network. In *WACV*, 2023.
- [29] 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.
- [30] Xiangtai Li, Wenwei Zhang, Jiangmiao Pang, Kai Chen, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. Video k-net: A simple, strong, and unified baseline for video segmentation. In *CVPR*, 2022.
- [31] Yanwei Li, Hengshuang Zhao, Xiaojuan Qi, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Fully convolutional networks for panoptic segmentation. In *CVPR*, 2021.
- [32] Huaijia Lin, Ruizheng Wu, Shu Liu, Jiangbo Lu, and Jiaya Jia. Video instance segmentation with a propose-reduce paradigm. In *ICCV*, 2021.
- [33] 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.
- [34] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun

- Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *arXiv preprint arXiv:2303.05499*, 2023.
- [35] Yang Liu, Idil Esen Zulfikar, Jonathon Luiten, Achal Dave, Deva Ramanan, Bastian Leibe, Aljoša Ošep, and Laura Leal-Taixé. Opening up open world tracking. In *CVPR*, 2022.
- [36] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [37] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. Hota: A higher order metric for evaluating multi-object tracking. *International journal of computer vision*, 129:548–578, 2021.
- [38] Jonathon Luiten, Idil Esen Zulfikar, and Bastian Leibe. Unovost: Unsupervised offline video object segmentation and tracking. In *WACV*, 2020.
- [39] Jiayu Miao, Xiaohan Wang, Yu Wu, Wei Li, Xu Zhang, Yunchao Wei, and Yi Yang. Large-scale video panoptic segmentation in the wild: A benchmark. In *CVPR*, 2022.
- [40] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *ICCV*, 2019.
- [41] 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.
- [42] Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan Yuille, Philip HS Torr, and Song Bai. Occluded video instance segmentation. *IJCV*, 2022.
- [43] Lu Qi, Jason Kuen, Yi Wang, Jiuxiang Gu, Hengshuang Zhao, Zhe Lin, Philip Torr, and Jiaya Jia. Open-world entity segmentation. In *arXiv preprint arXiv:2107.14228*, 2021.
- [44] Siyuan Qiao, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Vip-deeplab: Learning visual perception with depth-aware video panoptic segmentation. In *CVPR*, 2021.
- [45] Xuebin Qin, Hang Dai, Xiaobin Hu, Deng-Ping Fan, Ling Shao, and Luc Van Gool. Highly accurate dichotomous image segmentation. In *ECCV*, 2022.
- [46] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [47] Frano Rajiĉ, Lei Ke, Yu-Wing Tai, Chi-Keung Tang, Martin Danelljan, and Fisher Yu. Segment anything meets point tracking. In *arXiv preprint arXiv:2307.01197*, 2023.
- [48] Sucheng Ren, Wenxi Liu, Yongtuo Liu, Haoxin Chen, Guoqiang Han, and Shengfeng He. Reciprocal transformations for unsupervised video object segmentation. In *CVPR*, 2021.
- [49] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *ECCV*, 2020.
- [50] Jianping Shi, Qiong Yan, Li Xu, and Jiaya Jia. Hierarchical image saliency detection on extended cssd. In *TPAMI*, 2015.
- [51] Siyu Tang, Mykhaylo Andriluka, Bjoern Andres, and Bernt Schiele. Multiple people tracking by lifted multicut and person re-identification. In *CVPR*, 2017.
- [52] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020.
- [53] 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.
- [54] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *CVPR*, 2021.
- [55] Mark Weber, Jun Xie, Maxwell Collins, Yukun Zhu, Paul Voigtlaender, Hartwig Adam, Bradley Green, Andreas Geiger, Bastian Leibe, Daniel Cremers, et al. Step: Segmenting and tracking every pixel. In *NeurIPS*, 2021.
- [56] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation. In *CVPR*, 2022.
- [57] 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.
- [58] Bin Yan, Yi Jiang, Peize Sun, Dong Wang, Zehuan Yuan, Ping Luo, and Huchuan Lu. Towards grand unification of object tracking. In *ECCV*, 2022.
- [59] Bin Yan, Yi Jiang, Jiannan Wu, Dong Wang, Ping Luo, Zehuan Yuan, and Huchuan Lu. Universal instance perception as object discovery and retrieval. In *CVPR*, 2023.
- [60] Jinyu Yang, Mingqi Gao, Zhe Li, Shang Gao, Fangjing Wang, and Feng Zheng. Track anything: Segment anything meets videos. In *arXiv preprint arXiv:2304.11968*, 2023.
- [61] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *ICCV*, 2019.
- [62] Zongxin Yang, Yunchao Wei, and Yi Yang. Associating objects with transformers for video object segmentation. In *NeurIPS*, 2021.
- [63] Yi Zeng, Pingping Zhang, Jianming Zhang, Zhe Lin, and Huchuan Lu. Towards high-resolution salient object detection. In *ICCV*, 2019.