

---

# VITA: Video Instance Segmentation via Object Token Association

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce a novel paradigm for offline Video Instance Segmentation (VIS),  
2 based on the hypothesis that explicit object-oriented information can be a strong  
3 clue for understanding the context of the entire sequence. To this end, we pro-  
4 pose VITA, a simple structure built on top of an off-the-shelf Transformer-based  
5 image instance segmentation model. Specifically, we use an image detector as a  
6 means of distilling object-specific contexts into object tokens. VITA accomplishes  
7 video-level understanding by associating frame-level object tokens without using  
8 spatio-temporal backbone features. By effectively building relationships between  
9 objects using the condensed information, VITA achieves the state-of-the-art on  
10 VIS benchmarks with ResNet-50 backbone: 49.8 AP, 45.7 AP on YouTube-VIS  
11 2019 & 2021 and 19.6 AP on OVIS. Moreover, thanks to its object token-based  
12 structure that is disjoint from the backbone features, VITA shows several practical  
13 advantages that previous offline VIS methods have not explored - handling long and  
14 high-resolution videos with a common GPU and freezing a frame-level detector  
15 trained on image domain. Code will be made available.

## 16 1 Introduction

17 The goal of Video Instance Segmentation (VIS) is to predict both mask trajectories and categories  
18 of each object belonging to a set of predefined categories. Numerous studies have attained the  
19 goal in a variety of ways, but a notable innovation in terms of accuracy has been achieved by  
20 Transformer-based [25] architectures. Extending DETR [5] to the video domain, VisTR [26] made  
21 the first attempt to design an end-to-end model that jointly predicts object trajectories with their  
22 corresponding segmentation masks. By adopting this paradigm, subsequent studies [14, 28, 6, 32]  
23 also tackle the problem in a complete-offline manner: *video-in and video-out*.

24 The key message from the follow-up approaches [14, 28, 6, 32] is to effectively design core interac-  
25 tions between frames. In parallel with recent studies [10, 23, 35, 7, 20] that improve the accuracy in  
26 various tasks by localizing the attention scope of Transformer layers, the subsequent VIS methods  
27 suggest bounding the attention scope in the encoder [14, 32] or the decoder [28]. Specifically, they  
28 decompose the global attention by iteratively mixing two phases: intra-frame attention and inter-frame  
29 communication. Interestingly, the temporal interactions between frames are commonly achieved with  
30 only a small number of tokens, *e.g.*, memory tokens [14, 32], messenger tokens [32], and instance  
31 queries [28]. As a result, the question arises: “what information is important to understand a video?”

32 In this paper, we introduce Video Instance Segmentation via Object Token Association (VITA), a  
33 new offline VIS paradigm which suggests that a video can be effectively understood from a collection

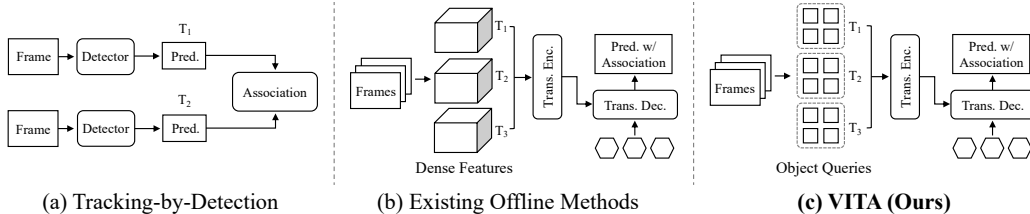


Figure 1: (a) Early-stage VIS methods divide the problem into two components, detection and association. (b) To alleviate the context-limited structure, complete-offline methods jointly track and segment instances in an end-to-end manner by employing dense spatio-temporal features. (c) On the other hand, our VITA is a new paradigm that directly leverages object queries for offline VIS.

34 of object-centric tokens. Existing offline methods [26, 14, 28, 6, 32] (Fig. 1 (b)) localize objects  
 35 in multiple frames by iteratively referring to dense spatio-temporal backbone features. However,  
 36 such methods show difficulties in handling long sequences as the myriad of dense reference features  
 37 hinders the Transformer layers from retrieving relevant information. With the motivation to devise  
 38 an effective method for the long-range understanding, we obtain clues from the traditional tracking-  
 39 by-detection paradigm (Fig. 1 (a)) and make two hypotheses: 1) an image object detector can fully  
 40 embody the context of an object into a feature vector (or a token); and 2) a video can be represented  
 41 by the relationship between the objects.

42 In this regard, VITA aims to parse an input video from the collection of object tokens without  
 43 the necessity of referencing dense spatio-temporal backbone (Fig. 1 (c)). Given the compactness  
 44 of the token representation, VITA can collect the object tokens over the whole video and directly  
 45 analyzes the collection using Transformer layers. This unique design enables the complete-offline  
 46 inference (i.e., video-in and video-out) even for extremely long videos. This also facilitates building  
 47 relationships between every detected object and successfully achieves global video understanding.  
 48 As a result, VITA achieves state-of-the-art performance on various VIS benchmarks.

49 We evaluate VITA on three popular VIS benchmarks, YouTube-VIS 2019 & 2021 [30] and OVIS [22].  
 50 With ResNet-50 [13] backbone, VITA achieves the new state-of-the-arts of 49.8 AP & 45.7 AP on  
 51 YouTube-VIS 2019 & 2021 and 19.6 AP on OVIS. Above all, VITA outperforms the previous best  
 52 approaches by 5.1 AP for YouTube-VIS 2021, which contains more complicated and long sequences  
 53 than YouTube-VIS 2019. VITA is the first offline method that presents the results on OVIS benchmark  
 54 that consists of long videos (the longest video has 292 frames) using a single 12GB GPU.

55 In addition to the performance, the design of VITA have several practical advantages over the previous  
 56 offline VIS methods. It can handle long and high-resolution videos so it does not require heuristics for  
 57 associating clip-level results. VITA can process 1392 frames at once regardless of video resolution  
 58 using a single 12GB GPU which is 11 times longer than IFC [14]. Moreover, VITA can be trained on  
 59 top of a parameter-frozen image detector without sacrificing the performance much. This property is  
 60 especially useful for the applications that cannot afford to store separated image and video instance  
 61 segmentation models. VITA takes only 6% additional parameters to extend the Swin-L detector.

## 62 2 Related Works

63 **Online VIS** approaches first predict individual tracklets within a local range window consisting of  
 64 a single or a few frames. After obtaining results from adjacent windows, they associate individual  
 65 tracklets of same identities by a hand-crafted or a learnable matching algorithm. MaskTrack R-  
 66 CNN [30] sets the groundwork for VIS research by proposing a simple tracking branch added  
 67 on a two-stage image instance segmentation model [12]. The methods [4, 31, 19] that follow the  
 68 tracking-by-detection paradigm (Fig. 1 (a)) measure the similarities between per-frame predictions,  
 69 then employ an association algorithm.

70 To deploy temporal context from multiple frames, [1, 2] design an architecture of predicting tracklets  
 71 within a local window and stitching the tracklets sequentially in a near-online manner. [15, 11] devise  
 72 a propagation paradigm that conjugates rich previous information stored in memories to facilitate

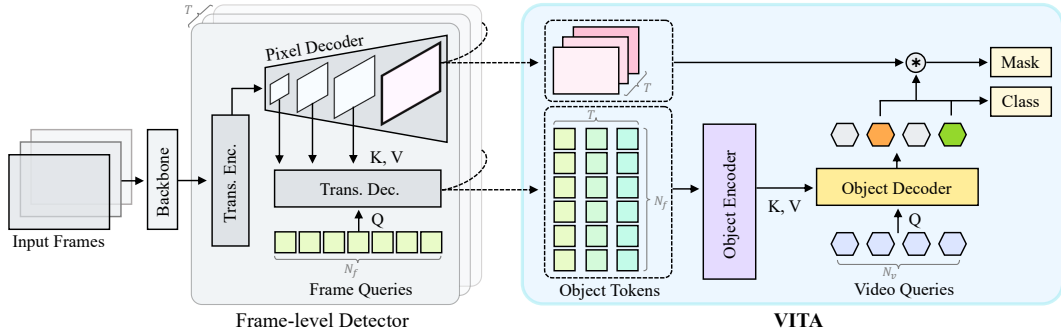


Figure 2: VITA takes only mask features and frame queries that are independently decoded by the frame-level detector for entire video sequence. By directly constructing temporal interactions between frame queries that encapsulate rich object-aware knowledge in spatial scenes, VITA yields mask trajectories with corresponding categories in an end-to-end manner.

73 online applications. EfficientVIS [27] introduces correspondence learning between adjacent tracklet  
 74 features and successfully runs in a cascaded manner which eliminates the hand-crafted tracklet  
 75 association.

76 **Offline VIS** architectures are proposed with the motivation of predicting mask trajectories through  
 77 a whole video sequence at once. VisTR [26] successfully extends DETR [5] to the VIS domain,  
 78 introducing a new paradigm of jointly tracking and segmenting instances. However, its dense self-  
 79 attention over the spatio-temporal inputs leads to explosive computations and memories. With the  
 80 motivation of relaxing the heavy computation of VisTR, IFC [14] adopts memory tokens to the  
 81 Transformer encoder and decodes clip-level object queries. By setting the frame-level encoder to be  
 82 independent and adopting the decoder of IFC, Mask2Former-VIS [6] records considerable perform-  
 83 ance on benchmarks by taking the advantage of its mask-oriented representation [7]. TeViT [32]  
 84 proposes a new backbone that efficiently exchanges temporal information internally based on Vision  
 85 Transformers [9] instead of the frame-wise CNN backbone. SeqFormer [28] decomposes the decoder  
 86 to be frame-independent, while building communication between different frames using instance  
 87 queries that are used for frame-wise detection. All these studies achieve promising performance by  
 88 referring to dense backbone features (Fig. 1 (b)). On the other hand, our VITA suggests a new offline  
 89 VIS paradigm that directly interprets a video from the collection of object tokens (Fig. 1 (c)).

90 **Global trackers** that aim to associate frame-level predictions across an entire sequence as a whole  
 91 are studied in the Multiple Object Tracking (MOT) community. Conventional approaches formulate  
 92 the problem as a graph optimization – interpreting each detection as a node and considering the edges  
 93 as possible connections between the nodes [33, 24, 3, 8]. Different from existing methods, GTR [34]  
 94 introduces a Transformer-based architecture that receives queries, then explicitly searches for the  
 95 predictions with the same identities.

### 96 3 Method

97 In this section, we first give a brief overview of Mask2Former [7], a frame-level detector for VITA.  
 98 Then, we introduce the architecture of our proposed VITA, which is built on top of Mask2Former.  
 99 Finally, we describe how VITA handles extremely long videos in a complete-offline manner.

#### 100 3.1 Frame-level Detector

101 In this paper, we adopt Mask2Former [7] for the frame-level detector which directly localizes instances  
 102 using masks without the necessity of bounding boxes. Following the set prediction mechanism of  
 103 DETR [5], the frame-level detector parse an input image  $H \times W$  using  $N_f$  object queries, which  
 104 we call *frame queries* ( $f \in \mathbb{R}^{C \times N_f}$ ) throughout this paper. Having the spatially encoded features  
 105 to be decoded by the frame queries through a Transformer decoder, each object in the image gets  
 106 represented as a  $C$ -dimensional vector. Then, the frame queries are used for both classifying and

107 segmenting their matched objects where the predictions are also used for auxiliary supervision for  
 108 VITA. Specifically, the frame-level detector generates two features for the frame-level predictions: 1)  
 109 dynamic  $1 \times 1$  convolutional weight from the frame queries; 2) per-pixel embeddings  $\mathcal{M} \in \mathbb{R}^{C \times \frac{H}{S} \times \frac{W}{S}}$   
 110 from the pixel decoder, where  $S$  is the stride of the feature map. Finally, the detector segments objects  
 111 by applying a simple dot product between the two embeddings.

### 112 3.2 VITA

113 We now propose the novel end-to-end video instance segmentation method VITA, which can be  
 114 largely divided into three phases (Fig. 2). First, VITA operates on top of the frame-level detector [7]  
 115 in a complete frame-independent manner; no inter-computation between frames is involved. Then,  
 116 the frame queries that hold object-centric information are collected throughout the whole video and  
 117 they embed video-level information by building communications between different frames using  
 118 Object Encoder. Finally, Object Decoder aggregates information from the frame queries to video  
 119 queries, which are eventually used for predicting categories and masks of objects in videos at once.

120 **Input of VITA.** Given an input video of  $T$  frames, the frame-level detector executes frame-by-frame  
 121 as previously explained. Among a number of intermediate embeddings that are generated by the  
 122 detector, the only features that are used by VITA are 1) frame queries  $\{f^t\}_{t=1}^T \in \mathbb{R}^{C \times T \times N_f}$  which  
 123 hold object-centric information; and 2) per-pixel embeddings  $\{\mathcal{M}^t\}_{t=1}^T \in \mathbb{R}^{C \times T \times \frac{H}{S} \times \frac{W}{S}}$  from the  
 124 pixel decoder.

125 **Object Encoder.** After the frame-wise detector dis-  
 126 tills the object-wise context into the frame queries,  
 127 Object Encoder aims to build temporal communi-  
 128 cation by employing self-attention along the tempo-  
 129 ral axis. First, Object Encoder gathers frame queries  
 130 from all frames and converts them to object tokens  
 131 through a linear layer. However, a naive self-attention  
 132 over the whole  $T N_f$  object tokens is not applica-  
 133 ble when processing long videos due to the quadra-  
 134 tic computational overhead of Transformers. Inspired  
 135 by Swin Transformer [20], we adopt window-based  
 136 self-attention layers that shift along the tempo-  
 137 ral dimension. As illustrated in Fig. 3, Object Encoder initially partitions object tokens  $\{f^t\}_{t=1}^T$  to the  
 138 temporal axis with local windows of size  $W$  without an overlap. By alternatively shifting the windows,  
 139 object tokens from different frames can exchange object-wise information which allows VITA to  
 140 both effectively and efficiently handle long sequences.

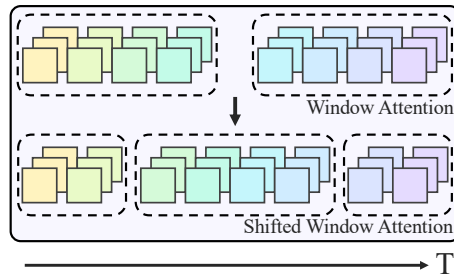


Figure 3: Illustration of an Object Encoder layer. Blocks with dashed line are local windows, and  $\square$  indicates an object token.

141 **Object Decoder and Output heads.** Two limitations of previous offline VIS methods [26, 14, 6]  
 142 are the ineffectiveness in handling dynamic scenes and the inability of processing long videos. For  
 143 example, although such methods obtain high accuracy when dealing with static and short videos  
 144 (YouTube-VIS 2019 [30]), they fail in either tracking objects or executing end-to-end on benchmarks  
 145 with dynamic and long videos (YouTube-VIS 2021 [30] and OVIS [22]). Both limitations are mainly  
 146 caused by the decoder, which parses object contexts directly from dense spatio-temporal features.  
 147 As recent studies [10, 23, 35] suggest, typical Transformer decoders show difficulties in retrieving  
 148 relevant information from global context. In the video domain, the number of backbone features  
 149 being referred to proportionally increases with the number of frames. Therefore, when handling  
 150 extremely long videos, the countless reference tokens result in both imprecise information retrieval  
 151 and intractable peak memories.

152 For the solution to the problem, we suggest Object Decoder which extracts information from the  
 153 object tokens, not the spatio-temporal backbone features. Implicitly embedding the context of objects,  
 154 object tokens can provide sufficient instance-specific information without the interference of dense  
 155 backbone features. Specifically, we employ  $N_v$  trainable video queries  $v \in \mathbb{R}^{C \times N_v}$  to decode object-

156 wise information from all object tokens  $\{f^t\}_{t=1}^T$  that are collected from all  $T$  frames. Receiving  
 157 much condensed input over naively taking dense spatio-temporal features, Object Decoder effectively  
 158 captures video contexts and aggregates relevant information into the video queries. As a result, Object  
 159 Decoder shows fast convergence speed while achieving high accuracy. Furthermore, the compact  
 160 input greatly saves memories, thus facilitates processing long and high-resolution videos.

161 From the decoded video queries  $v$ , VITA returns final predictions  $z = \{(p_i, m_i)\}_{i=1}^{N_v}$  using two  
 162 output heads similar to IFC [14]; the class head and the mask head. The class head is a single linear  
 163 classifier, which directly predicts class probabilities  $p \in \mathbb{R}^{N_v \times (K+1)}$  of each video query, where  
 164  $K + 1$  is the number of categories including an auxiliary label “no object” ( $\emptyset$ ). The mask head  
 165 dynamically generates mask embeddings  $w_v \in \mathbb{R}^{C \times N_v}$  per a video query, which corresponds to the  
 166 tracklet of an instance over all frames. Finally, the predicted mask logits  $m \in \mathbb{R}^{N_v \times T \times H \times W}$  can be  
 167 obtained from a matrix multiplication between  $w_v$  and  $\{\mathcal{M}^t\}_{t=1}^T$ .

### 168 3.3 Clip-wise losses

169 **Instance matching.** We search for optimal  
 170 pair indices between the predictions from  
 171 VITA and  $G_v$  ground-truth to remove post-  
 172 processing heuristics such as NMS. First, we  
 173 calculate costs from all possible pairs using  
 174 the cost function of Mask2Former [7] with a  
 175 simple extension of mask-related costs to the  
 176 temporal axis [14]. Then, from  $N_v \times G_v$  costs  
 177 of pairs, we follow DETR [5] and use Hun-  
 178 garian algorithm [16] for the optimal match-  
 179 ing as shown in Fig. 4 (b).

180 **Similarity loss.** Inspired by the initial VIS  
 181 approach (MaskTrack R-CNN [30]) where  
 182 the similarity loss is adopted to track in-  
 183 stances at different frames, we train video  
 184 queries and frame queries to be clustered in  
 185 the latent space by their identities. As shown  
 186 in Fig. 4 (a), our adopted frame-level detector [7] also searches for paired indices between  $N_f$   
 187 frame-wise predictions and  $G_f^t$  ground-truth objects at each  $t^{\text{th}}$  frame. The frame queries and  
 188 the video queries that are matched to ground-truths get collected and we embed the collection through a  
 189 linear layer. Then, we measure the similarity of all possible pairs using a simple matrix multiplication.  
 190 Finally, as shown in Fig. 4 (c), binary cross entropy is used to compute  $\mathcal{L}_{sim}$  between the predicted  
 191 similarities and the ground-truth where annotated to 1 for pairs of equal identities and 0 for vice-versa.

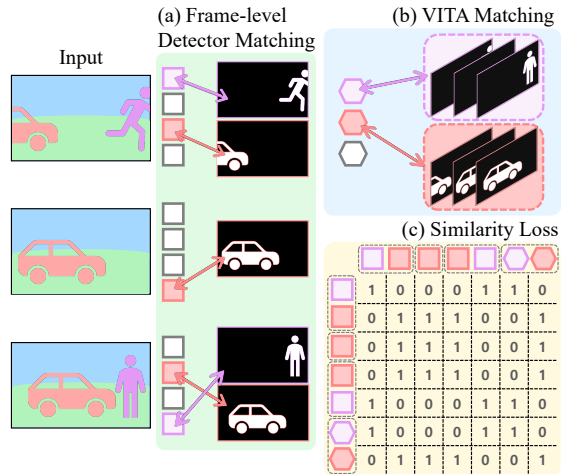


Figure 4: Similarity loss.  $\circ$  and  $\square$  indicate video query and frame query, respectively. Same color represents same GT instance ID.

192 **Total loss.** We attach the proposed module VITA on top of the frame-level detector, and all  
 193 components of the model get trained end-to-end. Note that not only video-level outputs from VITA  
 194 are used for the loss computation, but also per-frame outputs from the frame-level detector get  
 195 involved. Specifically, we use  $\mathcal{L}_f$  from [7] to calculate loss from the per-frame outputs to frame-wise  
 196 ground-truth. Extending the loss function of [7] to the temporal axis as similar to [14], we use outputs  
 197 from VITA  $z$  to calculate the video-level loss  $\mathcal{L}_v$ . Finally, we integrate all losses together as follows:  
 198  $\mathcal{L}_{total} = \lambda_v \mathcal{L}_v + \lambda_f \mathcal{L}_f + \lambda_{sim} \mathcal{L}_{sim}$ .

## 199 4 Experiments

### 200 4.1 Datasets

201 **YouTube-VIS 2019.** YouTube-VIS 2019 [30] is the first dataset proposed for VIS and contains 40  
 202 semantic categories. Mostly originated from Video Object Segmentation (VOS) datasets, the VIS

Table 1: Comparisons on YouTube-VIS 2019.

Method		Backbone [13]	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
(Near) Online	MaskTrack R-CNN [30]	ResNet-50	30.3	51.1	32.6	31.0	35.5
	MaskTrack R-CNN [30]	ResNet-101	31.8	53.0	33.6	33.2	37.6
	CrossVIS [31]	ResNet-50	36.3	56.8	38.9	35.6	40.7
	CrossVIS [31]	ResNet-101	36.6	57.3	39.7	36.0	42.0
	PCAN [15]	ResNet-50	36.1	54.9	39.4	36.3	41.6
	PCAN [15]	ResNet-101	37.6	57.2	41.3	37.2	43.9
	EfficientVIS [27]	ResNet-50	37.9	59.7	43.0	40.3	46.6
	EfficientVIS [27]	ResNet-101	39.8	61.8	44.7	42.1	49.8
	VISOLO [11]	ResNet-50	38.6	56.3	43.7	35.7	42.5
Offline	VisTR [26]	ResNet-50	35.6	56.8	37.0	35.2	40.2
	VisTR [26]	ResNet-101	38.6	61.3	42.3	37.6	44.2
	IFC [14]	ResNet-50	41.2	65.1	44.6	42.3	49.6
	IFC [14]	ResNet-101	42.6	66.6	46.3	43.5	51.4
	TeVIT [32]	MsgShifT	46.6	71.3	51.6	44.9	54.3
	SeqFormer [28]	ResNet-50	47.4	69.8	51.8	45.5	54.8
	SeqFormer [28]	ResNet-101	49.0	71.1	55.7	46.8	56.9
	SeqFormer [28]	Swin-L	59.3	82.1	66.4	51.7	64.4
	Mask2Former-VIS [6]	ResNet-50	46.4	68.0	50.0	-	-
	Mask2Former-VIS [6]	ResNet-101	49.2	72.8	54.2	-	-
	Mask2Former-VIS [6]	Swin-L	60.4	84.4	67.0	-	-
	VITA (Ours)	ResNet-50	49.8	72.6	54.5	49.4	61.0
		ResNet-101	51.9	75.4	57.0	49.6	59.1
Swin-L		63.0	86.9	67.9	56.3	68.1	

203 benchmark has a small number of unique instances (average 1.7 per video for the `train` set) and the  
204 categories of instances appearing in the same video are different in general. Also, the average length  
205 of videos in the `valid` set is short (27.4 frames), which enables existing complete-offline approaches  
206 to load a whole video and infer the benchmark at once.

207 **YouTube-VIS 2021.** In order to address more difficult scenarios, additional videos are included in  
208 YouTube-VIS2021 (794 videos for training and 129 videos for validation). In particular, a greater  
209 number of objects with confusing trajectories has been added (average 3.4 per video for the additional  
210 videos in the `train` set). However, the average length of the additional validation videos is still 39.7  
211 frames, which is not significantly increased compared to YouTube-VIS 2019.

212 **OVIS.** Under the same definition as YouTube-VIS, OVIS [22] specifically aims to tackle objects  
213 with heavy occlusions that are belonging to 25 semantic categories. In addition to the heavily occluded  
214 situation, OVIS has three challenging characteristics that are distinct from the YouTube-VIS datasets.  
215 First, although it has fewer categories than YouTube-VIS, much more instances appear in a single  
216 video (average 5.9 per video for the `train` set). Second, the instances with the same categories in the  
217 same video have almost similar appearances, thus approaches that rely heavily on visual cues often  
218 struggle to predict accurate trajectories. Finally, the average length of videos for the `valid` set is 62.7  
219 frames (the longest video has 292 frames) which is much longer than that of YouTube-VIS. Therefore,  
220 not only do previous approaches show relatively low accuracy, but all existing complete-offline VIS  
221 methods are not feasible to infer OVIS without hand-crafted association algorithms.

## 222 4.2 Implementation Details

223 Our method is implemented on top of `detectron2` [29]. All hyper-parameters regarding the frame-  
224 level detector are equal to the defaults of [7]. The total loss  $\mathcal{L}_{total}$  is balanced with  $\lambda_v$ ,  $\lambda_f$ , and  $\lambda_{sim}$   
225 where 1.0, 1.0, and 0.5, respectively. By default, Object Encoder is composed of three layers with the

Table 2: Comparisons with ResNet-50 backbone on YouTube-VIS 2021 and OVIS. † indicates using MsgShiT backbone.

Method	YouTube-VIS 2021					OVIS				
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
MaskTrack R-CNN [30]	28.6	48.9	29.6	26.5	33.8	10.8	25.3	8.5	7.9	14.9
CMaskTrack R-CNN [21]	-	-	-	-	-	15.4	33.9	13.1	9.3	20.0
STMASK [17]	31.1	50.4	33.5	26.9	35.6	15.4	33.8	12.5	8.9	21.3
CrossVIS [31]	34.2	54.4	37.9	30.4	38.2	14.9	32.7	12.1	10.3	19.8
IFC [14]	35.2	55.9	37.7	32.6	42.9	-	-	-	-	-
VISOLO [11]	36.9	54.7	40.2	30.6	40.9	15.3	31.0	13.8	11.1	21.7
TeVIT <sup>†</sup> [32]	37.9	61.2	42.1	35.1	44.6	17.4	34.9	15.0	11.2	21.8
SeqFormer [28]	40.5	62.4	43.7	36.1	48.1	-	-	-	-	-
Mask2Former-VIS [6]	40.6	60.9	41.8	-	-	-	-	-	-	-
<b>VITA (Ours)</b>	<b>45.7</b>	<b>67.4</b>	<b>49.5</b>	<b>40.9</b>	<b>53.6</b>	<b>19.6</b>	<b>41.2</b>	<b>17.4</b>	<b>11.7</b>	<b>26.0</b>

226 window size  $W = 6$ , and Object Decoder employs six layers with  $N_v = 100$  video queries. Having  
 227 VITA built on top of [7], we first train our model on the COCO [18] dataset following [7]. Then, we  
 228 train our method on the VIS datasets [30, 22] simultaneously with pseudo videos generated from  
 229 images [18] following the details of [28]. During inference, each frame is resized to a shorter edge  
 230 size of 360 and 480 pixels when using ResNet [13] and Swin [20] backbones, respectively. Note  
 231 that all reported scores in main results and ablation studies are the mean of five runs, and we use the  
 232 standard ResNet-50 [13] for the backbone unless specified.

### 233 4.3 Main Results

234 Using the popular VIS benchmarks – YouTube-VIS 2019 & 2021 [30] and OVIS [22] – we compare  
 235 VITA with state-of-the-art approaches following the standard evaluation metric [30].

236 **YouTube-VIS 2019.** Tab. 1 shows the comparison on YouTube-VIS 2019 dataset with backbones of  
 237 both CNN-based (ResNet-50 and 101 [13]) and Transformer-based (Swin-L [20]). Offline methods  
 238 can take two advantages over (near) online approaches: 1) they have a greater receptive field to  
 239 the temporal axis, and 2) they can avoid error propagation derived from hand-crafted association  
 240 algorithms. As a result, the tendency of offline methods with higher accuracy is clearly shown in the  
 241 table. Among the competitive offline models, our VITA sets a new state-of-the-art of 49.8 AP and  
 242 51.7 AP using CNN backbones, ResNet-50 and ResNet-101 respectively. In addition, with Swin-L  
 243 backbone, VITA achieves 63.0 AP outperforming all existing VIS methods.

244 **YouTube-VIS 2021.** We compare VITA with state-of-the-art methods on YouTube-VIS 2021  
 245 benchmark in Tab. 2. Note that the longest video in the valid set has 84 frames, thus previous  
 246 offline methods [14, 28, 6] can infer videos at once with GPUs with large memories. Above all, VITA  
 247 achieves the highest accuracy, 45.7 AP, which outperforms the previous state-of-the-art approach [6]  
 248 with a huge margin of 5.1 AP. Considering the accuracy gap on YouTube-VIS 2019, the results  
 249 demonstrate that VITA can effectively handle tricky scenarios, *e.g.*, numerous unique instances with  
 250 confusing trajectories. We hypothesize that the object-oriented design of VITA is more effective than  
 251 typical dense Transformer decoders in addressing such challenging scenes.

252 **OVIS.** In Tab. 2, we demonstrate the competitiveness of VITA on the challenging OVIS benchmark.  
 253 Due to the considerable lengths of videos – the longest video has 292 frames – existing offline  
 254 approaches [26, 14, 32, 28, 6] cannot process OVIS benchmark in their original design: video-in and  
 255 video-out. To the best of our knowledge, VITA is the first complete-offline approach to evaluate on  
 256 OVIS valid set. Thanks to its object token-based structure which is disjoint from backbone features,  
 257 VITA can process the benchmark *without* any hand-crafted association algorithm. Moreover, VITA  
 258 sets a new state-of-the-art performance of 19.6 AP, demonstrating the potential of the complete-offline  
 259 pipeline in long and complicated scenes.

Table 3: Impact of local windows of varying sizes in Object Encoder.

$W$	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
3	49.4	72.2	54.4	48.6	60.9
6	49.8	72.6	54.5	49.4	61.0
12	50.0	73.0	54.7	49.0	60.8
All	50.1	72.4	54.7	49.0	60.6

Table 4: Maximum number of frames that can be processed at once using a single Titan XP.

Method	Max Frames	
	360 × 640	720 × 1280
VisTR [26]	46	12
IFC [14]	123	38
Mask2Former-VIS [6]	81	20
VITA	$W = 3$	2677
(Ours)	$W = 6$	1392
	$W = 12$	741

#### 260 4.4 Ablation Studies

261 We provide a series of ablation studies using a ResNet-50 [13] backbone. All experiments are  
 262 conducted on YouTube-VIS 2019 [30] valid set except for Tab. 5 with OVIS [22] valid set.

263 **Attention window size.** Tab. 3 shows the performance of VITA with varying sizes of shifted  
 264 attention window  $W$  in Object Encoder during inference. The larger the window, the greater the  
 265 receptive field for the temporal axis in Object Encoder. The results suggest that larger window  
 266 sizes utilize information from multiple frames, which helps Object Encoder understand the context  
 267 of objects in videos. We set  $W = 6$  considering a trade-off between performance and inference  
 268 scalability.

269 **Maximum number of frames.** In Tab. 4, we calculate the maximum number of frames that VITA  
 270 can handle with respect to the various window sizes  $W$ , and compare it with existing complete-offline  
 271 VIS methods. To take into account the general environment, all results are computed using a single  
 272 12GB Titan XP GPU. As shown in results, existing methods have limitations in processing long  
 273 videos in a *video-in and video-out* manner. Clearly, the bottleneck of VisTR [26] is the encoder, where  
 274 the full spatio-temporal self-attention leads to a tremendous memory usage. IFC [14] alleviates the  
 275 computation of VisTR [26], achieving a higher number of input frames. However, IFC makes use of a  
 276 typical Transformer decoder that visits all dense spatio-temporal features. Therefore, IFC cannot infer  
 277 the OVIS [22] benchmark at once which contains a video of 292 frames. The problem gets aggravated  
 278 in Mask2Former-VIS [6] as the scope of the decoder is extended to multiple feature levels [7]. On  
 279 the other hand, VITA presents considerable frame numbers that can be inferred completely offline.  
 280 Furthermore, VITA is independent from input frame resolutions as each frame gets summarized into  
 281 compact object tokens. With input resolution of  $360 \times 640$  and  $W = 6$ , the maximum length of  
 282 sequence that VITA is able to process in complete-offline is about  $11 \times$  longer than IFC [14].

283 **Heuristic clip association.** Tab. 5 shows the results on  
 284 OVIS valid set of splitting a video into shorter clips  
 285 and associating clip-wise predictions through heuristic  
 286 matching. The length of the clip is set to be less than the  
 287 average length of videos of OVIS valid set (62.7). Then,  
 288 we associate outputs from different clips using mask IoU  
 289 score as the matching cost. We test with two matching  
 290 algorithms: Greedy and Hungarian. As shown in Tab. 5,  
 291 VITA demonstrates the best performance on the complete-  
 292 offline inference that use all the video frames at once.

Table 5: Use of different heuristic association algorithms on OVIS valid set.

Length	Algorithm	AP	AP <sub>50</sub>	AP <sub>75</sub>
36	Greedy	18.8	39.4	17.1
	Hungarian	18.4	38.9	16.3
48	Greedy	18.8	39.0	17.1
	Hungarian	19.1	39.1	17.4
All	None	19.6	41.2	17.4

293 **Convergence speed and Similarity loss.** Fig. 5 validates our claim of the faster convergence  
 294 speed and the effectiveness of the proposed Similarity loss. For a fair comparison, we report  
 295 average scores and standard deviations of five runs, each trained without pseudo videos, same as  
 296 Mask2Former-VIS [6]. Thanks to its object-centric design, VITA shows faster convergence than  
 297 Mask2Former-VIS. Furthermore, the use of Similarity loss leads to an additional accuracy gain of

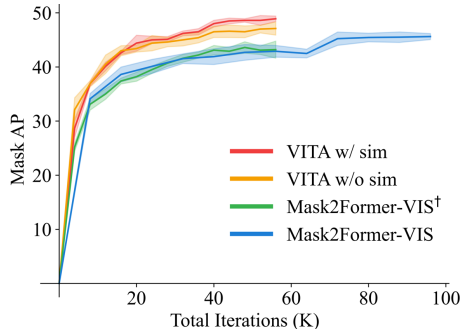


Figure 5: Train speed comparison with Mask2Former-VIS [6]. † indicates the same training setup with VITA.

Table 6: Freezing detector pretrained on COCO.

Backbone	Freeze	AP	AP <sub>50</sub>	AP <sub>75</sub>	AR <sub>1</sub>	AR <sub>10</sub>
ResNet-50	✓	49.8	72.6	54.5	49.4	61.0
		40.9	61.9	44.6	43.1	53.1
ResNet-101	✓	51.9	75.4	57.0	49.6	59.1
		43.2	64.4	48.7	46.1	55.9
Swin-L	✓	63.0	86.9	67.9	56.3	68.1
		53.4	75.9	58.7	51.9	64.3

298 1.8 AP. The results demonstrate that the loss mitigates the discrepancies between the embeddings of  
 299 equal identities, leading to better performance.

300 **Frozen frame-level detector.** In Tab. 6, we demonstrate the performance of VITA where the  
 301 frame-level detector is completely frozen. Specifically, while VITA gets trained on YouTube-VIS  
 302 2019, the frame-level detector [7] does not get updated from pretrained weights on COCO [18].  
 303 Note that among 40 categories in YouTube-VIS 2019 dataset, only 20 categories overlap with the  
 304 categories of COCO. Interestingly, though the frame-level detector remains completely frozen, VITA  
 305 achieves compelling results with various backbones. As shown in Tab. 1 and Tab. 6, VITA presents  
 306 a huge practicality as it surpasses all online approaches on top of the ResNet-50 backbone. This  
 307 strategy can be beneficial in various scenarios: 1) when the accuracy of image instance segmentation  
 308 should be kept while extending the network to the video domain, and 2) when having limited time  
 309 and GPUs to train models. The strategy can be especially useful in mobile applications that have  
 310 scarce storage for keeping two separate network parameters for image instance segmentation and  
 311 video instance segmentation. With additional 6% parameters, VITA successfully extends the frozen  
 312 Swin-L based frame-level detector to the video domain and it achieves great accuracy.

## 313 5 Conclusion

314 In this paper, we proposed VITA for offline Video Instance Segmentation. VITA is a simple model  
 315 built on top of the off-the-shelf image instance segmentation model [7]. Unlike existing offline  
 316 methods, VITA directly leverages object queries decoded by independent frame-level detectors.  
 317 We demonstrated that deploying object-oriented information is not only effective in improving  
 318 performance, but also has robust practicality for processing long and high-resolution videos - setting  
 319 state-of-the-art on popular VIS benchmarks, *e.g.*, YouTubeVIS-2019 & 2021 and OVIS. Moreover,  
 320 since VITA is designed to absorb spatial knowledge purely from image detectors, it shows fast  
 321 convergence and demonstrates competitive performance even if trained on frozen detectors. For the  
 322 ultimate long video understanding, we believe that devising an explicit design that takes temporal  
 323 information will be a promising future research direction.

## 324 Broader Impact

325 VITA is designed for the VIS task and focuses on processing long and high-resolution videos in an  
 326 end-to-end manner while achieving the state-of-the-art performance. We hope that VITA can have a  
 327 positive impact on many industrial areas such as video editing applications. We would like to note  
 328 that research on VIS must be aware of potential misuse that violates personal privacy.

329 **COCO [18], YouTube-VIS [30], OVIS [22], detectron2 [29] license:** Attribution 4.0 Interna-  
 330 tional, CC BY 4.0, CC BY-NC-SA 4.0, Apache-2.0

331 **References**

- 332 [1] Ali Athar, Sabarinath Mahadevan, Aljoša Ošep, Laura Leal-Taixé, and Bastianan Leibe. Stem-  
333 seg: Spatio-temporal embeddings for instance segmentation in videos. In *ECCV*, 2020.
- 334 [2] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances  
335 in video with mask propagation. In *CVPR*, 2020.
- 336 [3] Guillem Brasó and Laura Leal-Taixé. Learning a neural solver for multiple object tracking. In  
337 *CVPR*, 2020.
- 338 [4] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang,  
339 and Ling Shao. Sipmask: Spatial information preservation for fast image and video instance  
340 segmentation. In *ECCV*, 2020.
- 341 [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and  
342 Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020.
- 343 [6] Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexan-  
344 der G Schwing. Mask2former for video instance segmentation. *arXiv preprint arXiv:2112.10764*,  
345 2021.
- 346 [7] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar.  
347 Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022.
- 348 [8] Peng Dai, Renliang Weng, Wongun Choi, Changshui Zhang, Zhangping He, and Wei Ding.  
349 Learning a proposal classifier for multiple object tracking. In *CVPR*, 2021.
- 350 [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
351 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
352 An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- 353 [10] Peng Gao, Minghang Zheng, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fast convergence  
354 of detr with spatially modulated co-attention. In *ICCV*, 2021.
- 355 [11] Su Ho Han, Sukjun Hwang, Seoung Wug Oh, Yeonchool Park, Hyunwoo Kim, Min-Jung  
356 Kim, and Seon Joo Kim. Visolo: Grid-based space-time aggregation for efficient online video  
357 instance segmentation. In *CVPR*, 2022.
- 358 [12] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017.
- 359 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
360 recognition. In *CVPR*, 2016.
- 361 [14] Sukjun Hwang, Miran Heo, Seoung Wug Oh, and Seon Joo Kim. Video instance segmentation  
362 using inter-frame communication transformers. In *NeurIPS*, 2021.
- 363 [15] Lei Ke, Xia Li, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Prototypical  
364 cross-attention networks for multiple object tracking and segmentation. In *NeurIPS*, 2021.
- 365 [16] Harold W Kuhn. The hungarian method for the assignment problem. In *Naval research logistics*  
366 *quarterly*, 1955.
- 367 [17] Minghan Li, Shuai Li, Lida Li, and Lei Zhang. Spatial feature calibration and temporal fusion  
368 for effective one-stage video instance segmentation. In *CVPR*, 2021.
- 369 [18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,  
370 Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*,  
371 2014.

- 372 [19] Dongfang Liu, Yiming Cui, Wenbo Tan, and Yingjie Chen. Sg-net: Spatial granularity network  
373 for one-stage video instance segmentation. In *CVPR*, 2021.
- 374 [20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
375 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- 376 [21] Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan  
377 Yuille, Philip HS Torr, and Song Bai. Occluded video instance segmentation. *arXiv preprint*  
378 *arXiv:2102.01558*, 2021.
- 379 [22] Jiyang Qi, Yan Gao, Yao Hu, Xinggang Wang, Xiaoyu Liu, Xiang Bai, Serge Belongie, Alan  
380 Yuille, Philip HS Torr, and Song Bai. Occluded video instance segmentation: Dataset and iccv  
381 2021 challenge. In *NeurIPS Track on Datasets and Benchmarks*, 2021.
- 382 [23] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris M Kitani. Rethinking transformer-based  
383 set prediction for object detection. In *ICCV*, 2021.
- 384 [24] Siyu Tang, Mykhaylo Andriluka, Bjoern Andres, and Bernt Schiele. Multiple people tracking  
385 by lifted multicut and person re-identification. In *CVPR*, 2017.
- 386 [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
387 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 388 [26] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and  
389 Huaxia Xia. End-to-end video instance segmentation with transformers. In *CVPR*, 2020.
- 390 [27] Jialian Wu, Sudhir Yarram, Hui Liang, Tian Lan, Junsong Yuan, Jayan Eledath, and Gerard  
391 Medioni. Efficient video instance segmentation via tracklet query and proposal. In *CVPR*, 2022.
- 392 [28] Junfeng Wu, Yi Jiang, Wenqing Zhang, Xiang Bai, and Song Bai. Seqformer: a frustratingly  
393 simple model for video instance segmentation. *arXiv preprint arXiv:2112.08275*, 2021.
- 394 [29] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2.  
395 <https://github.com/facebookresearch/detectron2>, 2019.
- 396 [30] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *ICCV*, 2019.
- 397 [31] Shusheng Yang, Yuxin Fang, Xinggang Wang, Yu Li, Chen Fang, Ying Shan, Bin Feng, and  
398 Wenyu Liu. Crossover learning for fast online video instance segmentation. In *ICCV*, 2021.
- 399 [32] Shusheng Yang, Xinggang Wang, Yu Li, Yuxin Fang, Jiemin Fang, Wenyu Liu, Xun Zhao, and  
400 Ying Shan. Temporally efficient vision transformer for video instance segmentation. In *CVPR*,  
401 2022.
- 402 [33] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking  
403 using network flows. In *CVPR*, 2008.
- 404 [34] Xingyi Zhou, Tianwei Yin, Vladlen Koltun, and Phillip Krähenbühl. Global tracking transform-  
405 ers. In *CVPR*, 2022.
- 406 [35] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:  
407 Deformable transformers for end-to-end object detection. In *ICLR*, 2021.

408 **Checklist**

- 409 1. For all authors...
- 410 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
411 contributions and scope? [Yes]
- 412 (b) Did you describe the limitations of your work? [Yes] We describe the limitations of  
413 our work as a future research direction in Sec. 5.
- 414 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We  
415 included potential misuses in Broader Impact.
- 416 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
417 them? [Yes]
- 418 2. If you are including theoretical results...
- 419 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 420 (b) Did you include complete proofs of all theoretical results? [N/A]
- 421 3. If you ran experiments...
- 422 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
423 imental results (either in the supplemental material or as a URL)? [Yes] We listed  
424 the data and instructions in Sec. 4 and Appendix, and the code will be released upon  
425 acceptance.
- 426 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
427 were chosen)? [Yes] We specified details in Sec. 4.2.
- 428 (c) Did you report error bars (e.g., with respect to the random seed after running ex-  
429 periments multiple times)? [Yes] We reported error bars and score of five runs in  
430 Sec. 4.
- 431 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
432 of GPUs, internal cluster, or cloud provider)? [Yes] We included the type of resources  
433 used for training and inference in Appendix.
- 434 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 435 (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the papers  
436 and repositories that are used.
- 437 (b) Did you mention the license of the assets? [Yes] We mentioned the license at the end.
- 438 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
439 New assets are not included.
- 440 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
441 using/curating? [No]
- 442 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
443 information or offensive content? [No]
- 444 5. If you used crowdsourcing or conducted research with human subjects...
- 445 (a) Did you include the full text of instructions given to participants and screenshots, if  
446 applicable? [N/A] Not the scope of this paper.
- 447 (b) Did you describe any potential participant risks, with links to Institutional Review  
448 Board (IRB) approvals, if applicable? [N/A] Not the scope of this paper.
- 449 (c) Did you include the estimated hourly wage paid to participants and the total amount  
450 spent on participant compensation? [N/A] Not the scope of this paper.