

Online Video Instance Segmentation via Robust Context Fusion

Xiang Li, Jinglu Wang, Xiaohao Xu, Bhiksha Raj, *Fellow, IEEE* Yan Lu, *Senior Member, IEEE*

Video instance segmentation (VIS) aims at classifying, segmenting and tracking object instances in video sequences. Recent transformer-based neural networks have demonstrated their powerful capability of modeling spatio-temporal correlations for the VIS task. Relying on video- or clip-level input, they suffer from high latency and computational cost. We propose a robust context fusion network to tackle VIS in an online fashion, which predicts instance segmentation frame-by-frame with a few preceding frames. To acquire the precise and temporal-consistent prediction for each frame efficiently, the key idea is to fuse effective and compact context from reference frames into the target frame. Considering the different effects of reference and target frames on the target prediction, we first summarize contextual features through importance-aware compression. A transformer encoder is adopted to fuse the compressed context. Then, we leverage an order-preserving instance embedding to convey the identity-aware information and correspond the identities to predicted instance masks. We demonstrate that our robust fusion network achieves the best performance among existing online VIS methods and is even better than previously published clip-level methods on the Youtube-VIS 2019 and 2021 benchmarks.

In addition, visual objects often have acoustic signatures that are naturally synchronized with them in audio-bearing video recordings. By leveraging the flexibility of our context fusion network on multi-modal data, we further investigate the influence of audios on the video-dense prediction task, which has never been discussed in existing works. We build up an Audio-Visual Instance Segmentation dataset, and demonstrate that acoustic signals in the wild scenarios could benefit the VIS task.

Index Terms—video instance segmentation, multimodal learning

I. INTRODUCTION

The recently introduced video instance segmentation (VIS) problem receives increasing attention because of growing interest among researchers in the multimedia community. VIS aims at simultaneously classifying, segmenting and tracking objects in video sequences.

With strong capabilities for modeling long-range data dependencies, some transformer-based methods [55], [23] achieve impressive results. However, they deal with the input at the video- or clip- level, thus incurring high latency. The cost of modeling full spatio-temporal correlations is prohibitive for practical applications. In this work, we focus on *online* VIS for streaming applications. The problem setting is as follows: given a small set of preceding *reference* frames, the goal is to segment, classify and track object instances in each *target* frame.

In addition to the challenges of segmentation and classification in the image domain, online VIS should also address the problem of finding correspondences and fusing contexts in adjacent frames. Yang et al. [66] approach this problem by performing frame-based predictions independently and fusing the evidence across frames in a post-processing stage using sophisticated matching algorithms. Matching as post-processing has a high cost, and the final result may suffer from flickering because of neglecting feature-level correlations

across frames. Several subsequent methods [29], [2], [18], [48] fuse inter-frame features in the encoding stage. In particular, some methods [18], [29], [48] crop out ROI features to fuse context within local regions, but cropped features are isolated from the global context. To involve the global context, some methods [18], [60], [26] generally concatenate reference and target feature maps together; such “fusion” procedures do not distinguish reference and target frames, and thus the specific importance of the target frame is ignored. Moreover, as reference frames are usually similar to the target, spatio-temporal correlations among them could be highly redundant. Some reference features unrelated to the target could even mislead the target prediction.

Besides the contextual information contained in reference video frames, information in other modalities, such as audio, can also serve as reference context to guide the VIS task. Visual objects often have acoustic signatures that are naturally synchronized with them in audio-bearing video recordings. While audio-visual synchrony has been exploited in many contexts, most audio-visual representation learning methods handle signals in constrained environments, with clearly-evident audio-visual correspondences, such as music performance or lip reading. We deal with videos with audios in the wild where the correlation between audio and visual signals is difficult to establish. Previous methods [70], [53], [50], [45] have already demonstrated the correlation between audio and video modalities through sound localization and audio-visual separation, but there is no investigation on dense prediction done jointly with audio and visual modalities, potentially offering lost opportunities.

In this paper, we propose a robust context fusion method for online VIS. To preserve the location information of instances, we build our model on the crop-free transformer-based network. As shown in Figure 1, we employ an attention-based context fusion module for multi-modal contextual in-

Part of this work was done when Xiang Li was an intern at Microsoft Research Asia. Xiang Li is currently with Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213 USA. (email: xl6@andrew.cmu.edu)

Jinglu Wang and Yan Lu are with Microsoft Research Asia, Beijing, China. (email: jinglwa@microsoft.com; yanlu@microsoft.com)

This work was done when Xiaohao Xu was an intern at Microsoft Research Asia. Xiaohao Xu is currently with Department of Mechanical Engineering, Huazhong University of Science and Technology, Wuhan, China. (email: xxh11102019@outlook.com)

Bhiksha Raj is currently with School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213 USA. (email: bhiksha@cs.cmu.edu)

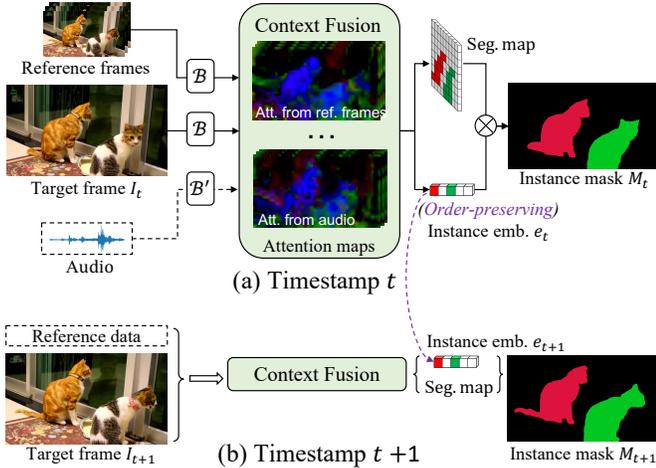


Fig. 1. We propose a robust context fusion (RCF) network for the online VIS task. Considering the redundancy and noise in the reference visual or audio context, the RCF module extracts compact and representative contexts from features (extracted by backbones \mathcal{B} and \mathcal{B}'), and then fuses them into the target feature with expressive attentions. The fused context is decoded into an instance code and segmentation maps. We directly compose the final instance masks by leveraging the *order-preserving* instance code and tracking instance identities without additional matching.

formation interaction, in which we compress the reference features to mitigate redundancy, improve robustness, and also introduce audio cues. An order-preserving instance code is further learned by the transformer decoder with a fixed-length instance query. Also, by leveraging the Lipschitz continuity of the network, we employ a matching-free instance identity tracking approach. The instance identity (shown as red and green colors in Figure 1) can be directly corresponded to the index of the instance code, thus greatly reducing the matching cost in the inference stage. From experimental and mathematical perspectives, we demonstrate that the order-preserving property of instance codes is tight under natural scenes even without any supervision. This can shed light on new directions of instance tracking. Our context fusion module is flexible to fuse multi-modal signals. We construct an unconstrained audio-visual dataset for the VIS task, and demonstrate that audio is helpful for the VIS task but the improvement is not significant due to the weak correlation between audio and visual signals in the wild scenarios. Our contributions are three-fold.

- A robust context fusion module for modeling compact and effective spatio-temporal correlations for online VIS. Our method achieves the state-of-the-art performance among online VIS methods on Youtube-VIS 2019 and 2021 benchmarks.
- A matching-free and supervision-free instance identity tracking method, and its corresponding mathematical explanation.
- We are the first to investigate audio effects on the video dense prediction task and contribute a corresponding dataset with synchronized audio-visual signals.

A preliminary version of our method has been published in [30]. In this manuscript, we make the following improvements.

- 1) We analyze the mathematical intuition behind the order-preserving property of instance code and conduct more experiments to understand its behavior.
- 2) We simplify the instance code generation to improve the network generalization, and at the same time the computational cost is reduced.
- 3) We extend the framework to accept audio-visual inputs, and present an effective method to fuse the multimodal context information leveraging the compatibility of the network.
- 4) We provide extensive quantitative and qualitative experimental results to show the performance of our framework in different settings and examine the effectiveness of the key modules in ablation studies.

II. RELATED WORK

Video instance segmentation. Video instance segmentation [66], [7], [38], [54], [54], [68], [31] requires classifying and segmenting each instance in a frame and assigning the same instance with the same identity across frames. There are mainly two types of methods for VIS tasks: online and offline.

Online VIS: Given a small set of preceding *reference* frames, online VIS aims to segment, classify and track object instances in each *target* frame. Mask-Track-RCNN [66] is the first attempt to address the VIS problem in an online setting. It extends the Mask-RCNN [20] with a tracking head to associate instance identities. SipMask [7] and SG-Net [38] build the tracking head on top of the modified one-stage still-image instance segmentation method FCOS [54] and BlenderMask [9], and achieve better speed and performance compared to MaskTrack-RCNN. CrossVIS [68] introduces the cross-over learning scheme and global instance embedding to learn better features for robust instance tracking and segmentation. HITF proposes inter-frame attention and intra-frame attention layers which bridges instance information across frames by an instance code.

Offline VIS: The offline methods handle the sequence-to-sequence prediction, where all frames are considered to be equivalent. Given a clip of video, offline VIS is to segment, classify and track object instances in all clip-level frames at the same time. VISTR [60] utilizes a transformer encoder on top of the convolutional backbone and matches instance identities by transformer decoder. IFC introduces a frame memory to reduce the computational cost for temporal aggregation, which decouples the segmentation in each frame and only communicates via frame memory.

Video object segmentation. Video object segmentation (VOS) aims to segment object masks across frames in a class-agnostic manner. Typically, additional cues are given to specify the target object. Semi-supervised VOS [67], [24], [61], [6], [46] gives the first frame object mask to specify the target object. RANet [61] proposes a ranking attention module to filter our irrelevant features based on the pixel-level similarity obtained from an encoder-decoder network. Space-Time-Memory networks (STM) [46] introduces a new paradigm that builds a memory bank for each object in the video and segments following objects by matching them to the memory bank. More recently, referring video object segmentation (R-VOS) emerges and attracts more and more attention because of its

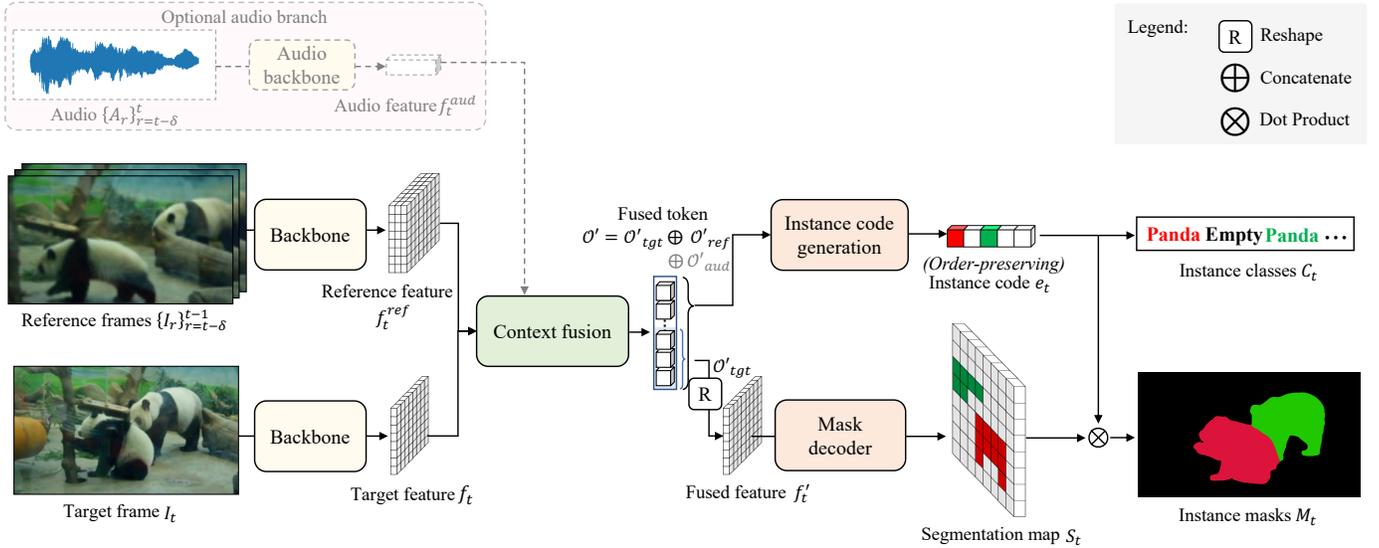


Fig. 2. **Overview of the proposed network at timestamp t .** For each target frame I_t , we consider reference frames $\{I_r\}_{r=t-\delta}^{t-1}$ as context for predicting the instance masks M_t and classes C_t . Both target feature f_t and reference feature f_t^{ref} are extracted by a shared backbone and fused in the context fusion module. The mask decoder decodes the fused target token \mathcal{O}'_{tgt} into the segmentation map S_t . The transformer decoder decodes the fused token \mathcal{O}' into an order-preserving instance embedding e_t in which each slot corresponds to a specific instance in I_t . The final predictions, instance classes C_t and instance masks M_t , are obtained from the instance embedding e_t and segmentation map S_t with a simple dot product operation. For the instance identity matching, we constrain the slot indices (indicated in red and green) of instance embedding e_t to represent the instance identity. The optional audio signal serves as another context and can be fused in the same way as reference frames.

strong ability to facilitate human-computer interaction. R-VOS [62], [5], [34], [51], [15], [32] aims to segment object masks throughout the entire video by giving a linguistic expression. URVOS first segments object masks by visual cues then select the referred one by referring expression. MTTR follows this paradigm while enable the multimodal fusion by a transformer encoder to enhance the semantic consensus between linguistic and visual modalities. ReferFormer directly segments the referred object by employing a language-conditioned instance query in the transformer decoder which avoids segmenting irrelevant objects.

Image Instance Segmentation Most of image instance segmentation methods adopt either bottom-up [12], [64], [58], [59], [57], [69], [10], [56] or top-down [20], [8], [39], [22] paradigm. For bottom-up methods, several representations can be adopted to represent instance identities, such as object centers [12], object-specific coefficients [8]. PolarMask [64] directly models instance contour by using 36 uniformly-spaced rays in polar coordinates, which can be assumed as a generalization of box representation that models each instance contour by 4 rays in polar coordinates. SOLO [58], [59] discriminates each instance by its location and size. Panoptic-Deeplab [12] models semantic segmentation, object center heatmap and center offset separately and then assembles those components to instance masks. SipMask [8] follows the previous one-stage method FCOS [54] to represent each instance by object-specific coefficients and corresponding mask prototypes. To achieve more accurate segmentation result, SipMask divides the input image to four parts and predicts masks in each part separately. For top-down methods, instance identity is mainly represented by the bounding-box of detected object. Mask R-CNN [20] employs a region proposal generation network

(RPN) equipped with RoIAlign feature pooling strategy and a feature pyramid networks (FPN) [36] to obtain fixed-sized features of each proposal. The pooled features are further used for bounding box prediction and mask segmentation. Followed by Mask R-CNN, several methods are proposed to improve pooling and confidence scoring strategy [33], [39], [22].

Audio-visual representative learning. Audio-visual representative learning aims to correspond the sound to their sources in the video frames. Sound localization [70], [53], [50], [65], [44], [14] is one of the more extensively explored tasks in this domain, which locates the sound sources of the audio recording in the image. Sound-of-pixel [70] conducts sound source localization and separation simultaneously on an instrument dataset, by estimating time-frequency masks on the audio spectrogram from the video, and recovering the separated sounds through an inverse Short Time Fourier Transform (STFT). A2V [53] leverages the attention mechanism to link the audio and video in unconstrained videos. Similarly, Senocak et al. [50] further expand the audio-visual localization task into semi-supervised scenarios and managed to improve the localization performance by introducing unsupervised loss. Audio-visual sound source (speaker) separation [1], [16], [42], [45], [19], [70] aims to isolate the target sound in a noisy scene. Recent methods leverage visual cues to guide the separation. Gabbay et al. [19] first feed the video frames into a video-to-speech model to predict the speaker’s voice from facial movements, then use it to separate the target speaker from sound mixtures.

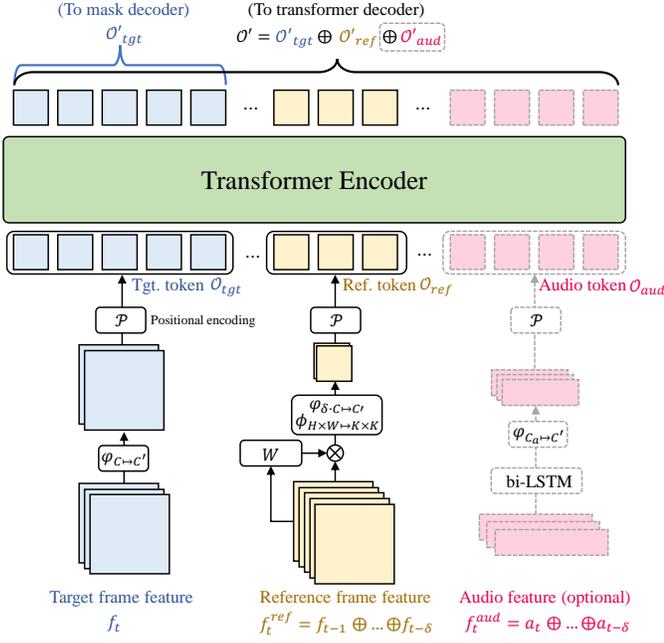


Fig. 3. **Robust context fusion module.** The target frame feature f_t is projected to a lower dimension, and then flattened and added with positional encoding to be tokens O_{tgt} . Considering the importance to the target prediction, features from reference frames f_t^{ref} are first reweighted by a learned mask W , and then compressed to lower spatial and channel dimensions as O_{ref} . The optional audio features are projected to a lower dimension as O_{aud} after bi-LSTM. We concatenate the source tokens and feed them to the transformer encoder to model their correlations.

III. METHOD

A. Overview

Pipeline overview. We first introduce our transformer-based network targeting online video instance segmentation, and then extend it by adding optional synchronized audio signals. The pipeline is illustrated in Figure 2. For each iteration t with a target frame I_t and reference frames $\{I_r\}_{r=t-\delta}^{t-1}$, we first extract the target feature f_t and reference features $\{f_r\}_{r=t-\delta}^{t-1}$ with a shared backbone. Considering the redundancy and noise in the pixel-wise correlations from reference to target frames, we compress the target and reference context as tokens O_{tgt} and O_{ref} respectively. Then, a transformer encoder in the robust context fusion (RCF) module fuses the original tokens as $O' = O'_{tgt} \oplus O'_{ref}$, where O'_{tgt} and O'_{ref} are fused tokens correspond to target and reference contexts. Afterwards, O'_{tgt} is decoded into the target segmentation map S_t and the overall O' is decoded into an instance code e_t , which represents order-preserving instance identities. We directly compose the final instance segmentation map M_t by dynamic convolution between S_t and e_t without any sophisticated matching algorithm. We elaborate our online VIS framework in Section III-B.

Our pipeline is flexible to aggregate multi-modal contexts. Audio is verified to be helpful for object recognition and localization, and we investigate the audio-visual correspondence and leverage it to improve the dense prediction (instance segmentation in this work). We align and fuse audio features f_t^{aud} in the RCF the same way as reference video frames.

The online audio-visual instance segmentation framework is detailed in Section III-C.

B. Online Video Instance Segmentation

Since our method performs online inference, we process one target frame I_t at each iteration with additional reference frames $\{I_r\}_{r=t-\delta}^{t-1}$. The target features f_t are extracted by a shared backbone, while reference features $\{f_r\}_{r=t-\delta}^{t-1}$ are obtained from previous iterations. Note that feature extraction for each frame is performed only once in the online inference stage, thus saving processing time. The extracted video feature of both target frame and reference frames are sent to the robust context fusion module for further interaction.

1) Robust Context Fusion (RCF)

We fuse the context information using compact and representative visual tokens for target and reference frames by taking their importance into consideration. The structure of RCF is illustrated in Figure 3.

Target token. Since the target frame contains the most important information about spatial and semantic cues, we only compress the feature map along the channel dimension and preserve the spatial dimension.

$$O_{tgt} = \mathcal{P}(\varphi_{C \rightarrow C'}(f_t)), \quad (1)$$

where $\varphi_{C \rightarrow C'}$ is a 1×1 conv layer to project the feature map $f_t \in \mathbb{R}^{C \times H \times W}$ to a lower dimension $\mathbb{R}^{C' \times H \times W}$, \mathcal{P} denotes operations to flatten the feature and add it with positional encoding.

Reference token. Reference tokens are used to enhance target tokens according to their correlations with the target, while their own representation is less important. We employ compact and representative reference tokens to alleviate matching noise and enhance the importance of the target. Feature compression in both spatial and channel dimensions is applied.

$$O_{ref} = \mathcal{P}(\varphi_{\delta \cdot C \rightarrow C'}(\phi_{H \times W \rightarrow K \times K}(W \cdot f_t^{ref}))), \quad (2)$$

where $f_t^{ref} = f_{t-1} \oplus \dots \oplus f_{t-\delta}$ is a concatenated feature of all reference features, $\phi_{H \times W \rightarrow K \times K}$ is a pooling or depth-wise conv layer to compress the spatial dimension from $H \times W$ to $K \times K$, $\varphi_{\delta \cdot C \rightarrow C'}$ is a 1×1 conv layer to compress the channel dimension from $\delta \cdot C$ to C' , $W = \varphi_{\delta \cdot C \rightarrow 1}(f_{ref})$ is a learned pixel-wise weight map (visualized in Figure 4 (b)). We get the final reference tokens as $O_{ref} \in \mathbb{R}^{C' \times K \cdot K}$.

Token fusion. The target tokens O_{tgt} and reference tokens O_{ref} are concatenated and fused with a transformer encoder, and then we get the fused tokens $O' = O'_{tgt} \oplus O'_{ref} \in \mathbb{R}^{C' \times (H \cdot W + K \cdot K)}$. The fused target tokens $O'_{tgt} \in \mathbb{R}^{C' \times H \cdot W}$ are further reshaped and fed into a mask decoder to generate the segmentation map S_t . In addition, the whole fused tokens O' are fed into a transformer decoder to extract instance-specific information, which is discussed in Section III-B2.

RCF analysis. We further analyze our design of RCF from theoretical and experimental perspectives. Following the transformer structure [3], we adopt standard the scaled dot-product attention, i.e., $\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$, in the transformer encoder of RCF. A softmax function is applied to obtain the weights on the values. Given a much smaller

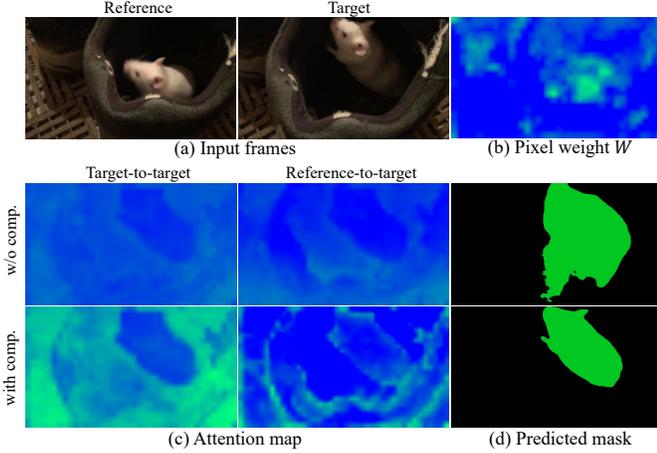


Fig. 4. **Visualization of intermediate components in RCF.** Given the hard cases (a) with confusing background, our design with token compression can generate more representative attention maps (c) and get better result (d).

token set (with size $|C' \times K \times K|$) of references, the target tokens ($|C' \times H \times W|$) dominate the attention in the transformer encoder. We thus get larger weights for target values of the compressed tokens than the uncompressed counterpart. The target-to-target attention map in Figure 4 (c) shows that target tokens in the transformer play more important roles than that without compression. Besides, the compressed tokens filter out irrelevant or noisy pixels and contain global information of references frames, thus are more discriminative and generate more representative attention. The reference-to-target attention is visualized in Figure 4 (c), where the compressed version looks more informative. According to the above two factors, our robust token compression can generate a better segmentation mask.

We underline that the goal of our online method is different from the offline methods [60], [23], and the network designs should be different accordingly. The offline methods tackle sequence-to-sequence prediction. The importance of all frames is equivalent, and compressing reference context is not mercenary. For the online setting, the target for each prediction is only the current frame, and importance-aware compression could benefit as discussed above.

2) Decoder

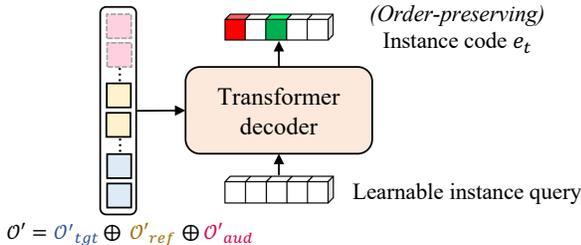


Fig. 5. **Illustration of instance code generation.** A learnable instance query is utilized to query the output of transformer encoder \mathcal{O}' to form the instance code e_t .

We generate instance code using a transformer decoder with a learnable instance query. As shown in Figure 5, the

learnable instance query is used to decode instance-specific information from the fused tokens \mathcal{O}' using a transformer decoder, where instance query is a set of learnable embedding as previous methods [60], [23], [63]. We term the decoded instance-specific features as instance code $e_t \in \mathbb{R}^{C_e \times N}$, where C_e and N are the channel dimension and length of the instance code respectively. For decoding the segmentation mask, we follow the FPN structure [36] to fuse the low-level features. Let the upsampled segmentation map be $S_t \in \mathbb{R}^{C_o \times H_o \times W_o}$ where C_o , H_o and W_o are the dimension, height and width of the upsampled segmentation map. After that, we leverage dynamic convolution [25] to obtain instance masks M_t from instance code. In particular, dynamic filters $\theta_t \in \mathbb{R}^{C_o \times N}$ are learned from instance code by two fully connected layers. The mask prediction $M_t \in \mathbb{R}^{N \times H_o \times W_o}$ can be computed as $M_t = \theta_t^T S_t$.

3) Loss Function

We assign each prediction with a ground-truth label then apply a set of loss functions between them. Given a set of predictions $\{\hat{p}_i(c), \hat{m}_i\}_{i=0}^N$ and a set of ground-truth $\{c_i, m_i\}_{i=0}^N$ (padded with \emptyset), we search for an assignment $\sigma \in \mathcal{S}_N$ with highest similarity, where $p_i(c)$ is the probability of class c (including \emptyset) and m_i is the mask of the i -th instance respectively. \mathcal{S}_N is a set of permutations of N elements. The similarity can be computed as

$$\text{Sim} = \mathbb{1}_{\{c_i \neq \emptyset\}} [\text{Dice}(m_{\sigma(i)}, \hat{m}_i) + \hat{p}_{\sigma(i)}(c_i)], \quad (3)$$

where $\mathbb{1}$ is an indicator function and Dice indicates the Dice loss [43]. The best assignment $\hat{\sigma}$ is solved by the Hungarian algorithm [28]. Given the best assignment $\hat{\sigma}$, the loss between ground-truth and predictions can be computed as

$$\mathcal{L} = \sum_{i=0}^N -\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \text{Dice}(\hat{m}_{\hat{\sigma}(i)}, m_i), \quad (4)$$

Following [13], we only consider the mask loss when the class prediction is not the empty class.

4) Instance Identity Matching

Unlike previous online methods [68], [7], [66] cropping out instances and using multiple cues (e.g., class, position and appearance) for matching instances across frames, we directly leverage the constraint that non-empty predictions from the same slot of the instance code have the same identity.

Order-preserving instance code. Our simple instance identity matching approach takes advantage of the Lipchitz continuity of our network (Appendix A) on the condition of bounded inputs. Given the Lipchitz continuity of our network $\Theta(\cdot)$ on the normalized image space \mathcal{I} , the relation of input and prediction can be represented as

$$0 \leq \|\Theta(I_t) - \Theta(I_{t-1})\|_p \leq \lambda \|I_t - I_{t-1}\|_p, \quad (5)$$

where λ is the Lipchitz constant and $\|\cdot\|_p$ represents p-norm with $p \in [1, \infty]$. In most cases, the discrepancy of adjacent frames is small, leading the discrepancy of their outputs $\Theta(\cdot)$ to be also small. Formally,

$$\lim_{\|I_t - I_{t-1}\|_p \rightarrow 0} \|\Theta(I_t) - \Theta(I_{t-1})\|_p = 0. \quad (6)$$

Since the order of output $\Theta(I_t) = \{\hat{p}_i(c), \hat{m}_i\}_{i=0}^N$ is directly linked to the order of the instance code e_t , if the order of e_t and e_{t-1} are different, the discrepancy $\|\Theta(I_t) - \Theta(I_{t-1})\|_p > \Delta$ could be large and does not satisfy Equation 6. Therefore, the order of instance code e_t would preserve given $\|I_t - I_{t-1}\|_p \rightarrow 0$. We utilize the preserved instance code order to track instance identities across frames. Different from previous offline methods [60], [23], [63], [30] that leverage losses to supervise the order-preserving, we claim that order-preserving is a property of our deep model and can work well in the online setting (with only adjacent frames) even without any supervision.

When $\|I_t - I_{t-1}\|_p \rightarrow 0$ holds, the order will preserve. While the this property of instance code is not constrained as tight as $\|I_t - I_{t-1}\|_p \rightarrow 0$ but related to the local smoothness of the network Θ . In other words, with a good local smoothness at a local region of Θ , even if $\|I_t - I_{t-1}\|_p \rightarrow 0$ does not hold, the order can still preserve. We provide with a relaxed deduction: If the order of the instance code changes, the Lipschitz continuity of the network will be violated. Let us denote $\mathcal{I}_{sub} \in \mathcal{I}$ as a local region containing input images I_t and I_{t-1} , where \mathcal{I} is the normalized input image space. Equation 5 (Lipschitz continuity) can be rewritten as:

$$\|\Theta(I_t) - \Theta(I_{t-1})\|_p < \lambda(\mathcal{I}_{sub}) \cdot \|I_t - I_{t-1}\|_p$$

where $\lambda(\mathcal{I}_{sub})$ is the Lipschitz constant of the network at \mathcal{I}_{sub} . For the **cases of order changes**, we denote $\epsilon_c = \|I_t - I_{t-1}\|_p$ as the input discrepancy, and $\Delta_c = \|\Theta(I_t) - \Theta(I_{t-1})\|_p$ as the output discrepancy. As the order changes, Δ_c can be very **large** according to the network design (Line 498-500). Then, the Lipschitz continuity takes the form:

$$\frac{\Delta_c}{\epsilon_c} < \lambda(\mathcal{I}_{sub})$$

If the network is well trained, Lipschitz constant $\lambda(\mathcal{I}_{sub})$ at a local region is always small to ensure the network smoothness. Since Δ_c is large and ϵ_c is a constant determined by the data, $\frac{\Delta_c}{\epsilon_c}$ can probably be larger than $\lambda(\mathcal{I}_{sub})$, which violates the Lipschitz continuity. Therefore, we can conclude that the order preserving property is mainly determined by $\lambda(\mathcal{I}_{sub})$ which reflects the local smoothness of the network, rather than purely by the input discrepancy $\|I_t - I_{t-1}\|_p$. In the case of $p = 1$, when the order changes, $\frac{\Delta_c}{\epsilon_c} \sim O(10^3)$ in Youtube-VIS dataset. Thus, the order can keep when the $\lambda(\mathcal{I}_{sub}) < O(10^3)$, which is not a tight constraint.

We empirically verify that the order-preserving matching solves most cases. However, there are still exceptional cases we cannot assume $\|I_t - I_{t-1}\|_p \rightarrow 0$ and $\|\Theta(I_t) - \Theta(I_{t-1})\|_p$ could also be small if tiny objects exist. Therefore, to enhance the model robustness, if the mask in i -th slot in time t has the IoU larger than 0.5 with mask in j -th slot in time $t - 1$, we directly assume they share the same identity without considering the order. Visualization analysis of fired slots in instance code will be discussed in Section IV-B1.

C. Online Audio-Visual Instance Segmentation

The video and audio data are naturally synchronized and contain the homogeneous semantic and spatial information of

the sound source. Compared to video data, audio stores the information in a more compact representation while 2D spatial locations of sound sources in the video frames can still be conveyed by audio cues [70], [53], [50], [45], [47].

VIS task requires semantic and spatial information of the objects which *de facto* corresponds to the information conveyed in the audio modality. We investigate the benefit of introducing audio data into VIS task in the following.

1) Audio-Visual Instance Segmentation Dataset

To investigate the influence of the audio data on the video-level fine-grained prediction tasks, we collect a novel dataset, AVIS, containing synchronized audio and video clips. The data is collected from publicly available videos with 20 vocal categories overlapped with Youtube-VIS dataset. AVIS contains 1427 videos with synchronized raw audio recordings, which are further randomly split by 75% and 25% for training and validation usage. All videos are annotated with high-quality instance-level labels in the format of Youtube-VIS dataset.

2) Audio Feature Extraction

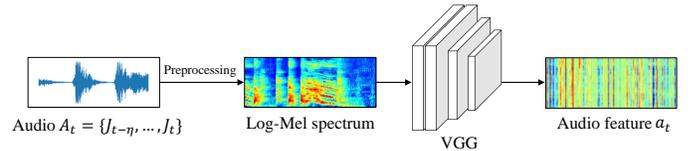


Fig. 6. **Audio processing.** The raw audio signal A_t is preprocessed by resampling, STFT, log mel transformation [40], and then fed into the pre-trained VGG [52] model (without the last linear layer) to extract feature a_t .

Since audio data has much higher sampling rate than video data, we combine multiple audio frames to one image. Given the image sequence $\{I_t\}_{t=1}^T$ and audio sequence $\{J_t\}_{t=1}^T$, we combine $A_t = \{J_{t-\eta}, \dots, J_t\}$ audio frames for each image I_t where η is the audio frame length for each image. We first resample A_t to 16 kHz mono then compute the spectrum using magnitudes of the Short-Time Fourier Transform (STFT) with a window size of 25 ms, a window hop of 10 ms, and a periodic Hann window. We pad the audio sequence to ensure the output has the same length as the input. The mel spectrogram is computed by mapping the spectrogram to 64 mel bins covering the range 125-7500 Hz. A stabilized log mel spectrogram is further computed by applying $\log(\text{mel-spectrum} + 0.01)$ where the offset is used to avoid taking a logarithm of zero. Figure 6 illustrates the audio feature extraction process. The generated log mel spectrogram [40] is fed to the pre-trained VGG [52] network to generate audio features. Note that we remove the last linear layer in the VGG to obtain high dimension features. The extracted audio feature for time t is denoted as a_t .

Audio features contain contextual information which relates to the position and semantic categories of the desired instances [70]. However, since audio understanding also depends on the contextual information in audio from prior frames, we introduce reference audio features as reference video features. Figure 2 shows an example of considering δ reference audio frames. The extracted features are sent to the context fusion module for cross-modal fusion.

Method	Pipeline	Backbone	FPS	Latency (s)	AP	AP50	AP75	AR1	AR10
VisTR [60]	Offline	ResNet-101	43.5	>1.9	38.6	61.3	42.3	37.6	44.2
MaskProp [4]	Offline	ResNet-101	-	-	42.5	-	45.6	-	-
IFC [23]	Offline	ResNet-101	89.4	6.4	44.6	69.2	49.5	44.0	52.1
Mask2Former [11]	Offline	ResNet-101	-	-	49.2	72.8	54.2	-	-
SeqFormer* [63]	Offline	ResNet-101	-	-	49.0	71.1	55.7	46.8	56.9
SeqFormer* [63]	Offline	Swin-L	-	-	59.3	82.1	66.4	51.7	64.4
MaskTrack R-CNN [66]	Online	ResNet-101	28.6	0.2	31.9	53.7	32.3	32.5	37.7
SipMask++[7]	Online	ResNet-101	27.8	0.2	35.0	56.1	35.2	36.0	41.2
SG-Net [38]	Online	ResNet-101	-	-	36.3	57.1	39.6	35.9	43.0
CrossVIS [68]	Online	ResNet-101	35.6	0.19	36.6	57.3	39.7	36.0	42.0
STMASK [29]	Online	ResNet-101	23.4	0.21	36.8	56.8	38.0	34.8	41.8
PCAN [26]	Online	ResNet-101	<15	>0.23	37.6	57.2	41.3	37.2	43.9
HITF [30]	Online	HITF	< 8	>0.29	41.3	61.5	43.5	42.7	47.8
Ours	Online	ResNet-101	23.9	0.21	40.8	63.4	43.5	42.4	46.7
Ours	Online	Swin-B	9.0	0.28	44.9	66.2	48.1	44.3	48.3
Ours	Online	Swin-L	5.3	0.36	47.6	70.3	50.8	46.1	52.3

TABLE I

COMPARISON TO THE STATE-OF-THE-ART OFFLINE AND ONLINE VIS METHODS ON THE YOUTUBE-VIS-2019 VALIDATION SET. THE INFERENCE SPEED IS MEASURED ON A SINGLE NVIDIA V100 GPU WITH $batchsize = 1$. * INDICATES THAT VALUES ARE IMPORTED FROM A PREPRINT.

3) Audio Context Fusion

We correspond the audio features to pixel-wise visual feature maps with cross-modal attention in the transformer encoder. Similar to tokens from visual features, we create audio tokens to support subsequent context fusion.

Audio token. We combine the overall audio features as the reference audio context $f_t^{aud} = a_{t-\delta} \oplus \dots \oplus a_t \in \mathbb{R}^{C_a \times (1+\delta)}$, where $C_a = 4096$. A two-layer bi-directional Long Short Time Memory (bi-LSTM) [21] network with hidden size 512 is leveraged to aggregate the temporal information. After that, we project the audio features f_t^{aud} to a lower dimension C' with two fully-connected layers and ReLU activation. Thus, we get the audio token $O_{aud} \in \mathbb{R}^{C' \times (1+\delta)}$.

IV. EXPERIMENT

In this section, we will elaborate on the dataset, implementation details and experiment results for our online VIS and AVIS frameworks.

A. Implementation Details.

Training. We implement our method in the PyTorch framework. Following previous methods [18], [35], we first pre-train our model with both Youtube-VIS and overlapped categories on MS-COCO dataset [37] then finetune the model on the Youtube-VIS dataset. We train our model for 60k iterations with a ‘‘poly’’ learning rate policy with the learning rate $(1 - \frac{iter}{iter_{max}})^{0.9}$ for each iteration with an initial learning rate of 0.0006 for all ResNet backbones and 0.0003 for all Swin backbones in experiments. We adopt $batchsize = 16$ and an AdamW [41] optimizer with $weight\ decay = 10^{-4}$ for all ResNet backbones and $weight\ decay = 10^{-2}$ for all Swin backbones is leveraged. A learning rate multiplier of 0.1 is applied to ResNet backbones, and 1.0 is applied to transformer backbones. Multi-scale training is adopted to obtain a strong baseline.

Inference. The image is resized to 640×360 during inference. To obtain the final object segmentation result, we first binarize

the $M_t \in \mathbb{R}^{N \times H \times W}$ by a threshold of 0.5. Then we filter out slots with a class probability less than 0.4 and keep the remaining ones as the predictions at time t . **Dataset.** We evaluate our method on the two extensively used VIS datasets, Youtube-VIS-2019 and Youtube-VIS-2021, as well as our newly constructed AVIS dataset.

- **Youtube-VIS-2019** has 40 categories, 4,883 unique video instances. There are 2,238 training videos, 302 validation videos, and 343 test videos in it.
- **Youtube-VIS-2021** is an improved version of Youtube-VIS-2019, which contains 8,171 unique video instances. There are 2,985 training videos, 421 validation videos, and 453 test videos in this dataset.
- **AVIS** has 20 overlapped categories with Youtube-VIS dataset, 2390 unique video instances. There are 1427 videos with synchronized raw audio recordings in it. To the best of our knowledge, AVIS is the first dataset with densely annotated instance-level masks matched to corresponding audio recordings.

Metrics. The evaluation metric for this task is defined as the area under the precision-recall curve with different IoUs as thresholds.

B. Video Instance Segmentation Results

We first present the main results and ablation experiments of VIS task on Youtube-VIS-2019 and Youtube-VIS-2021 dataset without audio involved.

1) Main results

We compare our method with state-of-the-art methods in this section.

Quantitative result. We compare our method against state-of-art VIS methods on **Youtube-VIS-2019** dataset in Table I. (1) Compared to online methods: Our method achieves the best performance of 40.8 mAP when using the same ResNet-101, outperforming the previous start-of-art method PCAN [26] by a large margin of 3.2 mAP. The recent method HITF [30]

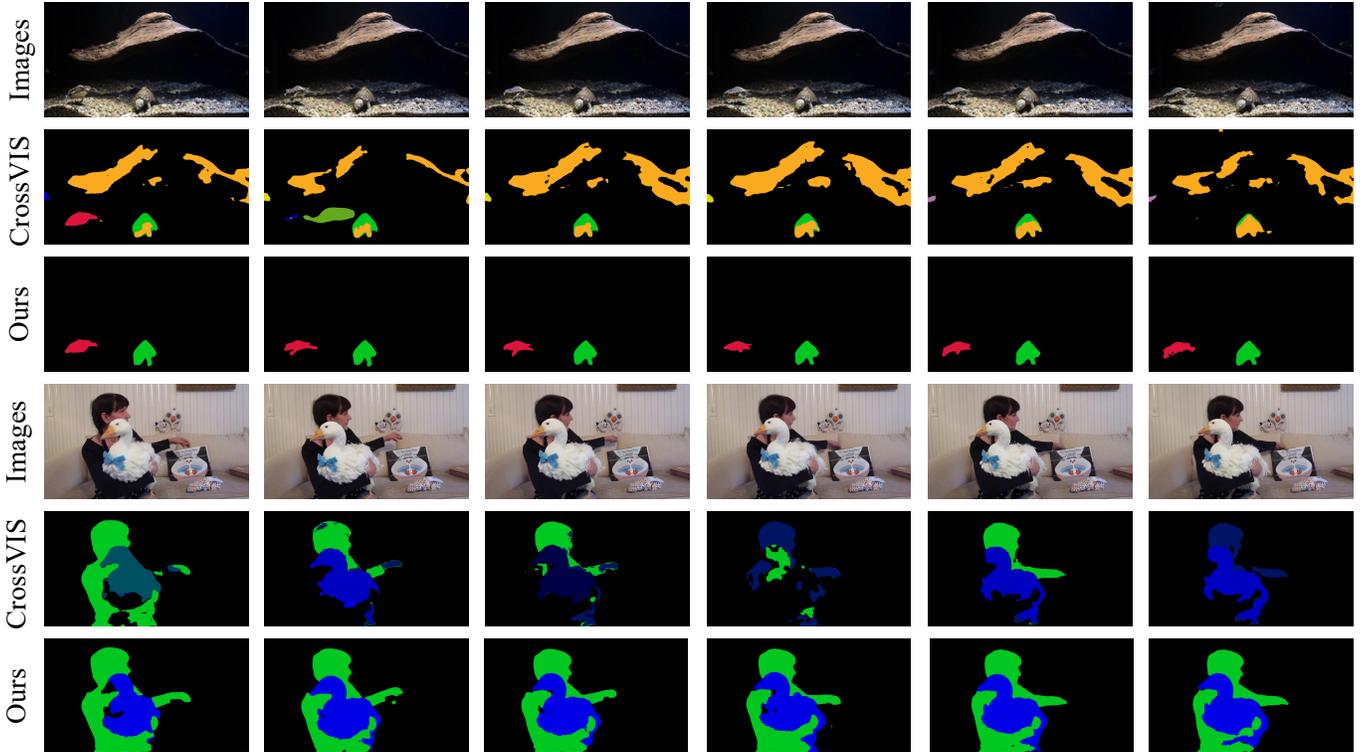


Fig. 7. Qualitative comparison to the state-of-the-art online video instance segmentation method CrossVIS [68] on the **Youtube-VIS-2019** validation set.

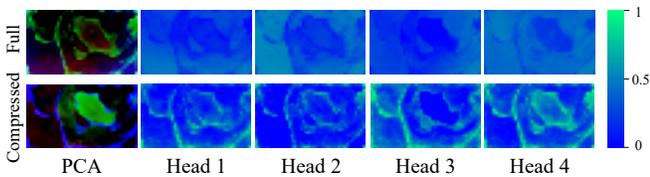


Fig. 8. **Visualization of reference-to-target attention map.** The attention maps from compressed tokens are more representative.

leverages a stronger backbone with a higher resolution input, we compare it with our results using the Swin-B backbone for fairness. Our method outperforms HITF by 3.6 mAP with a faster inference speed. (2) Compared to offline methods: Our method eclipses another transformer-based method VisTR [60] even if it takes video-level input. The results of SeqFormer [63] is imported from preprint, which uses a stronger deformable-transformer [71] to conduct temporal fusion, while other settings are similar to VisTR [60]. Although it achieves an impressive result, we consider it mainly due to the adoption of a stronger deformable transformer. We also compare our method against the state-of-the-art VIS methods on **Youtube-VIS 2021** in Table II. Since Youtube-VIS-2021 is a newly introduced dataset, there are only a few methods for comparison. Our method achieves the best performance among online methods.

Qualitative result. We compare the qualitative result of our method against CrossVIS [68] on Youtube-VIS-2019 val set

in Figure 7. The result indicates that CrossVIS fails to detect or track instances in occluded scenarios, while our methods successfully maintain robust high-quality segmentation and tracking performance. More qualitative results are shown in the supplementary material. As shown in Figure 8, we compare the reference-to-target (R2T) attention map in the transformer encoder using full and compressed reference tokens. The PCA is conducted among all 8 heads of R2T attention map and keeps three main components for RGB channels. We show the first 4 heads of the R2T attention. The attention maps of compressed tokens are sharper in the instance region than those of full tokens, which can also be verified by the PCA map.

Inference time. Online VIS is mainly for streaming applications. In the video streaming pipeline, the input video is received frame by frame, and latency (including video streaming and model processing time) is the essential time measurement. We further discuss the latency of both online and offline methods for fair comparison. Let $t_{stream} = 1/FPS_{stream}$ denote the streaming time of each frame and $t_{model} = 1/FPS_{model}$ denote the average model processing time for each frame. For online methods, the final latency of each frame is

$$Latency_{online} = t_{stream} + t_{model}$$

For offline methods, we need to wait till the N_f -th frame of the clip comes, where N_f is the frame number of the processed clip. The latency of the last frame (the whole clip) is

$$Latency_{offline} = t_{stream} * N_f + t_{model} * N_f$$

Method	Backbone	AP	AP50	AP75	AR1	AR10
MaskTrack	ResNet-50	28.6	48.9	29.6	26.5	33.8
SipMask-VIS	ResNet-50	31.7	52.5	34.0	30.8	37.8
CrossVIS	ResNet-50	34.2	54.4	37.9	30.4	38.2
CrossVIS	ResNet-101	35.2	56.3	36.4	33.9	40.6
HITF	HITF	35.8	56.3	39.1	33.6	40.3
Ours	ResNet-101	37.6	56.0	41.4	35.7	42.7

TABLE II

COMPARISON TO THE STATE-OF-THE-ART ONLINE METHODS ON **YOUTUBE-VIS-2021** VALIDATION SET.

Method	AP	AP50	AP75	AR1	AR10
full-feature Enc. [60]	35.6	58.8	37.3	36.6	41.3
full-feature Dec. [60]	36.6	57.7	38.0	36.9	41.8
reference-token	40.8	63.4	43.5	42.4	46.7

TABLE IV

COMPARISON OF DIFFERENT TEMPORAL FUSION METHODS ON **YOUTUBE-VIS-2019** VALIDATION SET.

Inst. Queries	AP	AP50	AP75	AR1	AR10
10	40.8	63.4	43.5	42.4	46.7
15	40.4	63.9	44.2	40.3	45.3
20	40.5	64.6	43.6	41.2	46.8
25	40.3	63.3	44.2	40.4	46.1

TABLE VI

IMPACT OF NUMBER OF INSTANCE QUERIES ON **YOUTUBE-VIS-2019** VALIDATION SET.

Weight W	AP	AP50	AP75	AR1	AR10
\times	40.1	63.2	42.1	40.3	46.3
\checkmark	40.8	63.4	43.5	42.4	46.7

TABLE VIII

IMPACT OF PIXEL-WISE WEIGHT W ON THE TOKEN FUSION.

Ref. Frame	AP	AP50	AP75	AR1	AR10
0	37.3	58.4	40.1	37.6	43.0
1	40.8	63.4	43.5	42.4	46.7
2	40.4	63.8	42.3	40.5	47.0
3	40.7	64.8	44.0	39.2	45.5
4	39.8	62.1	43.6	40.6	45.8
5	37.2	59.5	39.1	38.4	42.5

TABLE III

IMPACT OF THE REFERENCE FRAME NUMBER ON **YOUTUBE-VIS-2019** VALIDATION SET.

Token Size	AP	AP50	AP75	AR1	AR10
1	38.3	58.3	39.8	38.0	43.1
4	40.8	63.4	43.5	42.4	46.7
16	37.8	59.4	38.8	37.0	42.9
64	36.3	58.5	37.9	36.4	42.0

TABLE V

IMPACT OF THE REFERENCE-TOKEN SIZE ON **YOUTUBE-VIS-2019** VALIDATION SET.

Supervision	AP	AP50	AP75	AR1	AR10
\checkmark	40.5	63.5	42.6	40.5	46.7
\times	40.8	63.4	43.5	42.4	46.7

TABLE VII

IMPACT OF SUPERVISION ON THE ORDER-PRESERVING PROPERTY.

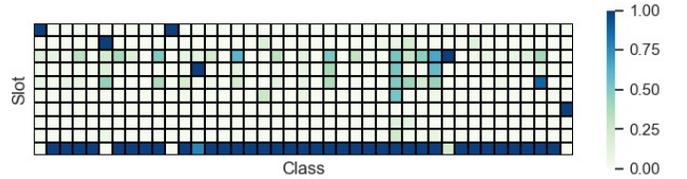


Fig. 9. **Illustration of fired slots and class correlation of instance code.** The heatmap is normalized by each class.

As the videos in Youtube-VIS dataset is of 6 FPS, the latency of our method is 0.171 s, while the latency of the offline method IFC with a clip containing $N_f = 36$ frames is 6.4 s ($FPS_{model} = 89.4$ reported in [13]).

As shown in Table I, we compare both FPS and latency of our method and previous methods. Our method achieves the best trade-off of accuracy and latency among online methods.

Analysis of fired slot and class correlation of instance code.

We consider the slot is fired it has a class probability larger than 0.4 and only keep predictions in fired slots as final output. As shown in Figure 9, we visualize the correlation of fired slots in instance code and each class. We noticed that most of objects are predicted from the same slot while, for several classes, e.g., person, duck and earless seal, they are predicted in stand-alone slots. Those classes are either most commonly appeared or challenging classes in the dataset. In this way, we consider the network encodes some class-specific information in the slot to improve instance segmentation quality.

2) Ablation Experiments

We conduct ablation studies on Youtube-VIS-2019 to show the effectiveness of different components of our method.

Fusion method. To investigate the importance of compressing reference features, we compare the performance of using reference tokens against using two baseline settings that directly fuse full reference features in transformer encoder and

transformer decoder respectively. To directly fuse features in transformer encoder, we compute reference tokens similarly to the target token. To fuse features in transformer decoder, we extract reference tokens \mathcal{O}'_{ref} by separate transformer encoders and concatenate them before the transformer decoder. As shown in Table IV, both baseline settings will result in a plummet in performance. Two reasons may account for the decrease. First, the instance position and appearance in the target frame differ from those in reference frames, which can mislead the network if giving reference frames the same importance as the target frame. Second, in transformer encoder, the attention of each layer is normalized by a softmax among all inputs. However, the mask prediction only needs fine-grained target frame information. Therefore, the softmax function downplays the importance of the target frame and introduces noises from reference frames.

Reference token size. The optimal reference token size is a trade-off between information compression loss and target importance gain. We conduct experiments on different token sizes. As shown in Table V, the token size of 4 achieves the best performance.

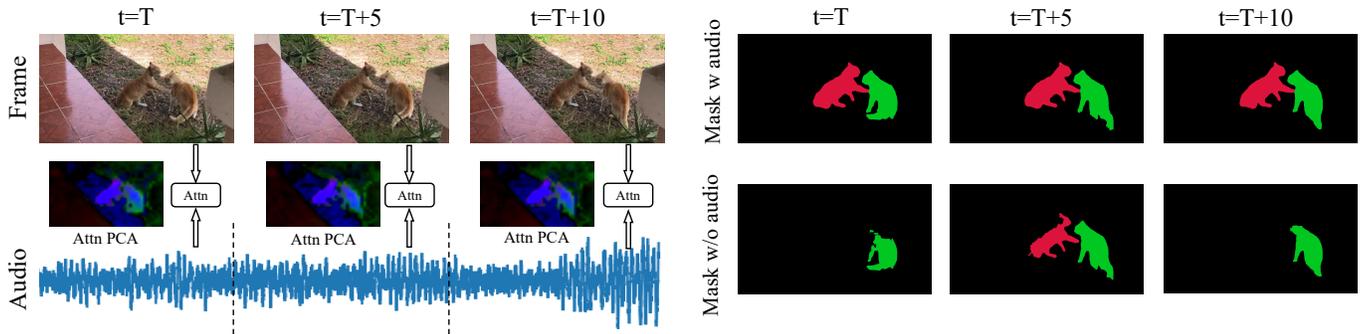


Fig. 10. Qualitative comparison of the performance with and w/o audio inputs. The multi-head attention map is compressed to three dimension by PCA for visualization.

Reference frame number. To investigate the importance of temporal information, we conduct an ablation study on training with different reference frames. As shown in Table III, as the reference frame number varies from 0 to 4, the mAP first increases to 40.8 then saturates and decreases. This is because the frames in Youtube-VIS dataset are provided at 6 FPS, thus introducing long-term references may bring in irrelevant information.

Number of instance queries. The number of instance queries indicates the maximum detected instances in a frame. However, there may be a varying number of instances in a frame within a video clip. As shown in Table VI, we notice that the model is robust to different instance query numbers even when they are severely redundant. The robustness to redundant queries enables the model to stay effective in scenarios that have extremely few instances.

Supervision on the order of instance code. Previous online VIS method HITF [30] explicitly gives supervision to keep the order-preserving. However, as we proved, the order-preserving is actually a property of the network. As shown in Table VII, we notice that the supervision of instance orders cannot improve the performance.

Pixel-wise weight W for the token fusion. Pixel-wise weight W aims to help the reference feature focus on the foreground area. As shown in Table VIII, we ablate on the influence of the pixel-wise weight W on the final segmentation performance. We notice that the performance drops 0.7 mAP if we disable the weight in the token fusion.

C. Audio-Visual Instance Segmentation Results

Method	Backbone	AP	AP50	AP75	AR1	AR10
CrossVIS	ResNet-50	42.2	59.1	47.4	42.1	49.9
CrossVIS	ResNet-101	44.7	60.4	50.5	42.0	49.7
Ours	ResNet-50	46.6	63.3	51.3	44.2	51.5
Ours	ResNet-101	48.6	63.9	52.5	44.6	52.9

TABLE IX

COMPARISON TO STATE-OF-THE-ART VIDEO INSTANCE SEGMENTATION ON AVIS VAL SET.

We report the performance of using ResNet-50 and ResNet-101 backbone on AVIS dataset. As shown in Table XI, we also compare the CrossVIS baseline which only leverages the

Ref. number	AP	AP50	AP	AP50	p-value
	w/o Audio-token		with Audio-token		
1	42.9	60.2	44.8 (+1.9)	60.3	0.30
2	44.3	60.9	46.0 (+1.7)	63.8	0.26
3	45.6	61.9	46.6 (+1.0)	63.3	0.65
4	45.6	62.2	45.7 (+0.1)	62.7	-
5	45.2	62.9	45.0 (-0.2)	62.6	-

TABLE X

IMPACT OF INTEGRATING AUDIO-TOKEN WITH DIFFERENT REFERENCE FRAMES ON AVIS VALIDATION SET.

video data. Our method outperforms CrossVIS baseline by 4.4 and 3.9 mAP with ResNet-50 and ResNet-101 backbone respectively.

To explore the benefit of introducing audio modality to VIS task, we conduct experiments on AVIS dataset to compare the models trained with audio inputs and without audio inputs. Figure 10 shows an example where audio improves the segmentation results. The attention between audio and video modality manages to locate the instance which is weakly distinguishable against the background. As shown in Table XII, we notice a slight gain from audio with a small reference frame number. However, as the reference frame number increases, the gain begins to saturate. There are several reasons accounting for the marginal gain. For example, there are some sound events that do not last for the whole video duration. Thus, the audio fusion is meaningless for several time steps. Moreover, since we add audio data serially with fixed reference audio frames, it is hard to construct a long-term correlation between audios thus the sound event cannot be fully utilized. To explore the statistical robustness of the gain, we conduct a t-test between corresponding experiments. As shown in Table XII, the p-value indicates the gains are not statistically significant. Therefore, we can safely conclude that while the audio may potentially benefit the video segmentation task, in the current setting, their effect in unconstrained videos is limited.

V. CONCLUSION

We propose a robust context fusion module for the online VIS task, which corresponds and fuses compact and effective reference features to the target features in a transformer encoder. We observe that the importance-aware compression

for reference context is critical in the online setting because the impact of frames are different for the target prediction. In addition, we leverage an order-preserving instance code to track instance identities, thus avoiding time-consuming matching algorithms. The mathematical explanation indicates that order-preserving is the natural property of the network and can work even without any explicit supervision, which may shed light on new directions of instance tracking. Our pipeline is flexible and permits multi-modal data. The benefits of using audio context to VIS are seen to be limited in the online setting due to weak correlation between audio and visual data in the wild scenes.

REFERENCES

- [1] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. The conversation: Deep audio-visual speech enhancement. *arXiv preprint arXiv:1804.04121*, 2018. 3
- [2] Ali Athar, Sabarinath Mahadevan, Aljosa Osep, Laura Leal-Taixé, and Bastian Leibe. Stem-seg: Spatio-temporal embeddings for instance segmentation in videos. In *European Conference on Computer Vision*, pages 158–177. Springer, 2020. 1
- [3] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3286–3295, 2019. 4
- [4] Gedas Bertasius and Lorenzo Torresani. Classifying, segmenting, and tracking object instances in video with mask propagation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9739–9748, 2020. 7
- [5] Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object segmentation with multimodal transformers. *arXiv preprint arXiv:2111.14821*, 2021. 3
- [6] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 221–230, 2017. 2
- [7] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. Sipmask: Spatial information preservation for fast image and video instance segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 1–18. Springer, 2020. 2, 5, 7
- [8] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. Sipmask: Spatial information preservation for fast instance segmentation, 2020. 3
- [9] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. Blendmask: Top-down meets bottom-up for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8573–8581, 2020. 2
- [10] Liang-Chieh Chen, Raphael Gontijo Lopes, Bowen Cheng, Maxwell D Collins, Ekin D Cubuk, Barret Zoph, Hartwig Adam, and Jonathon Shlens. Naive-student: Leveraging semi-supervised learning in video sequences for urban scene segmentation. In *European Conference on Computer Vision*, pages 695–714. Springer, 2020. 3
- [11] Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexander G Schwing. Mask2former for video instance segmentation. *arXiv preprint arXiv:2112.10764*, 2021. 7
- [12] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12475–12485, 2020. 3
- [13] Bowen Cheng, Alexander G Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *arXiv preprint arXiv:2107.06278*, 2021. 5, 14
- [14] Ying Cheng, Ruize Wang, Zhihao Pan, Rui Feng, and Yuejie Zhang. Look, listen, and attend: Co-attention network for self-supervised audio-visual representation learning. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3884–3892, 2020. 3
- [15] Zihan Ding, Tianrui Hui, Shaofei Huang, Si Liu, Xuan Luo, Junshi Huang, and Xiaoming Wei. Progressive multimodal interaction network for referring video object segmentation. *The 3rd Large-scale Video Object Segmentation Challenge*, page 7, 2021. 3
- [16] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. *arXiv preprint arXiv:1804.03619*, 2018. 3
- [17] Herbert Federer. *Geometric measure theory*. Springer, 2014. 13
- [18] Yang Fu, Linjie Yang, Ding Liu, Thomas S Huang, and Humphrey Shi. Compfeat: Comprehensive feature aggregation for video instance segmentation. *arXiv preprint arXiv:2012.03400*, 2020. 1, 7
- [19] Aviv Gabbay, Ariel Ephrat, Tavi Halperin, and Shmuel Peleg. Seeing through noise: Visually driven speaker separation and enhancement. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3051–3055. IEEE, 2018. 3
- [20] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2, 3
- [21] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 7
- [22] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring r-cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6409–6418, 2019. 3
- [23] Sukjun Hwang, Miran Heo, Seoung Wug Oh, and Seon Joo Kim. Video instance segmentation using inter-frame communication transformers. *arXiv preprint arXiv:2106.03299*, 2021. 1, 5, 6, 7
- [24] Suyog Dutt Jain, Bo Xiong, and Kristen Grauman. Fusionseg: Learning to combine motion and appearance for fully automatic segmentation of generic objects in videos. In *2017 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 2117–2126. IEEE, 2017. 2
- [25] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. *Advances in neural information processing systems*, 29:667–675, 2016. 5
- [26] Lei Ke, Xia Li, Martin Danelljan, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Prototypical cross-attention networks for multiple object tracking and segmentation. *arXiv preprint arXiv:2106.11958*, 2021. 1, 7
- [27] Hyunjik Kim, George Papamakarios, and Mnih. The lipschitz constant of self-attention. In *International Conference on Machine Learning*, pages 5562–5571. PMLR, 2021. 13
- [28] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 5
- [29] Minghan Li, Shuai Li, Lida Li, and Lei Zhang. Spatial feature calibration and temporal fusion for effective one-stage video instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11215–11224, 2021. 1, 7
- [30] Xiang Li, Jinglu Wang, Xiao Li, and Yan Lu. Hybrid instance-aware temporal fusion for online video instance segmentation. *arXiv preprint arXiv:2112.01695*, 2021. 2, 6, 7, 10
- [31] Xiang Li, Jinglu Wang, Xiao Li, and Yan Lu. Video instance segmentation by instance flow assembly. *arXiv preprint arXiv:2110.10599*, 2021. 2
- [32] Xiang Li, Jinglu Wang, Xiaohao Xu, Xiao Li, Yan Lu, and Bhiksha Raj. R2vos: Robust referring video object segmentation via relational multimodal cycle consistency. *arXiv e-prints*, pages arXiv:2207.2022.3, 2022. 3
- [33] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2359–2367, 2017. 3
- [34] Chen Liang, Yu Wu, Yawei Luo, and Yi Yang. Clawcranenet: Leveraging object-level relation for text-based video segmentation. *arXiv preprint arXiv:2103.10702*, 2021. 3
- [35] Huaijia Lin, Ruizheng Wu, Shu Liu, Jiangbo Lu, and Jiaya Jia. Video instance segmentation with a propose-reduce paradigm. *arXiv preprint arXiv:2103.13746*, 2021. 7
- [36] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 3, 5
- [37] 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 *European conference on computer vision*, pages 740–755. Springer, 2014. 7
- [38] Dongfang Liu, Yiming Cui, Wenbo Tan, and Yingjie Chen. Sg-net: Spatial granularity network for one-stage video instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9816–9825, 2021. 2, 7
- [39] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages

- 8759–8768, 2018. [3](#)
- [40] Beth Logan. Mel frequency cepstral coefficients for music modeling. In *International Symposium on Music Information Retrieval*. Citeseer, 2000. [6](#)
- [41] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. [7](#), [13](#)
- [42] Rui Lu, Zhiyao Duan, and Changshui Zhang. Listen and look: Audio-visual matching assisted speech source separation. *IEEE Signal Processing Letters*, 25(9):1315–1319, 2018. [3](#)
- [43] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. IEEE, 2016. [5](#)
- [44] Trisha Mittal, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. Emotions don’t lie: An audio-visual deepfake detection method using affective cues. In *Proceedings of the 28th ACM international conference on multimedia*, pages 2823–2832, 2020. [3](#)
- [45] Giovanni Morrone, Sonia Bergamaschi, Luca Pasa, Luciano Fadiga, Vadim Tikhonoff, and Leonardo Badino. Face landmark-based speaker-independent audio-visual speech enhancement in multi-talker environments. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6900–6904. IEEE, 2019. [1](#), [3](#), [6](#)
- [46] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9226–9235, 2019. [2](#)
- [47] Jie Pu, Yannis Panagakis, Stavros Petridis, and Maja Pantic. Audio-visual object localization and separation using low-rank and sparsity. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2901–2905. IEEE, 2017. [6](#)
- [48] 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. *arXiv preprint arXiv:2102.01558*, 2021. [1](#)
- [49] Ralph Schmidt. Multiple emitter location and signal parameter estimation. *IEEE transactions on antennas and propagation*, 34(3):276–280, 1986. [13](#)
- [50] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4358–4366, 2018. [1](#), [3](#), [6](#)
- [51] Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmentation network with a large-scale benchmark. In *European Conference on Computer Vision*, pages 208–223. Springer, 2020. [3](#)
- [52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [6](#)
- [53] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 247–263, 2018. [1](#), [3](#), [6](#)
- [54] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019. [2](#), [3](#)
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [1](#)
- [56] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. *arXiv preprint arXiv:2012.00759*, 2020. [3](#)
- [57] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *European Conference on Computer Vision*, pages 108–126. Springer, 2020. [3](#)
- [58] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *European Conference on Computer Vision*, pages 649–665. Springer, 2020. [3](#)
- [59] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic, faster and stronger. *arXiv preprint arXiv:2003.10152*, 2020. [3](#)
- [60] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8741–8750, 2021. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [9](#)
- [61] Ziqin Wang, Jun Xu, Li Liu, Fan Zhu, and Ling Shao. Ranet: Ranking attention network for fast video object segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3978–3987, 2019. [2](#)
- [62] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation. *arXiv preprint arXiv:2201.00487*, 2022. [3](#)
- [63] Junfeng Wu, Yi Jiang, Wenqing Zhang, Xiang Bai, and Song Bai. Seqformer: a frustratingly simple model for video instance segmentation. *arXiv preprint arXiv:2112.08275*, 2021. [5](#), [6](#), [7](#), [8](#)
- [64] Enze Xie, Peize Sun, Xiaoog Song, Wenhai Wang, Xuebo Liu, Ding Liang, Chunhua Shen, and Ping Luo. Polarmask: Single shot instance segmentation with polar representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12193–12202, 2020. [3](#)
- [65] Haoming Xu, Runhao Zeng, Qingyao Wu, Mingkui Tan, and Chuang Gan. Cross-modal relation-aware networks for audio-visual event localization. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3893–3901, 2020. [3](#)
- [66] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5188–5197, 2019. [1](#), [2](#), [5](#), [7](#)
- [67] Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos K Katsaggelos. Efficient video object segmentation via network modulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6499–6507, 2018. [2](#)
- [68] Shusheng Yang, Yuxin Fang, Xinggang Wang, Yu Li, Chen Fang, Ying Shan, Bin Feng, and Wenyu Liu. Crossover learning for fast online video instance segmentation. *arXiv preprint arXiv:2104.05970*, 2021. [2](#), [5](#), [7](#), [8](#)
- [69] Tien-Ju Yang, Maxwell D Collins, Yukun Zhu, Jyh-Jing Hwang, Ting Liu, Xiao Zhang, Vivienne Sze, George Papandreou, and Liang-Chieh Chen. Deeplab: Single-shot image parser. *arXiv preprint arXiv:1902.05093*, 2019. [3](#)
- [70] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *Proceedings of the European conference on computer vision (ECCV)*, pages 570–586, 2018. [1](#), [3](#), [6](#)
- [71] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [8](#)

VI. LIPSCHITZ CONTINUITY OF OUR NETWORK

We first define the notion of Lipschitz continuity, and proceed to the Lipschitz continuity of fully-connected, convolutional and self-attention layers. We refer to the proof in [27] with bounded inputs.

a) *High-level Proof:* Continuous functions are Lipschitz continuous given bounded inputs w.r.t. p -norm where $p \in [1, \infty]$.

Definition 1: Given two metric spaces $(\mathcal{X}, d_{\mathcal{X}})$ and $(\mathcal{Y}, d_{\mathcal{Y}})$, a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is called Lipschitz continuous (or K -Lipschitz) if there exists a constant $K \geq 0$ such that

$$d_{\mathcal{Y}}(f(x), f(x')) \leq K d_{\mathcal{X}}(x, x') \quad (7)$$

The smallest such K is the Lipschitz constant of f , denoted $\text{Lip}(f)$.

In the following proof, the $d_{\mathcal{X}}$ and $d_{\mathcal{Y}}$ are induced by a p -norm $\|x\|_p := (\sum_i |x_i|^p)^{1/p}$. Following [27], we emphasise the dependence of the Lipschitz constant on the choice of p -norm by denoting it as $\text{Lip}_p(f)$. In this case, it follows directly from Definition 1 that the Lipschitz constant is given by

$$\text{Lip}_p(f) = \sup_{x \neq x' \in \mathbb{R}^n} \frac{\|f(x) - f(x')\|_p}{\|x - x'\|_p} \quad (8)$$

A. Fully-connected/convolutional Layers

We describe how Lipschitz constant of fully-connected networks (FC) and convolutional neural networks (CNN) are bounded, using the fact that both are compositions of linear maps and pointwise non-linearities.

Theorem 1: ([17]) Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be differentiable and Lipschitz continuous under a choice of p -norm $\|\cdot\|_p$. Let $J_f(x)$ denote its total derivative (Jacobian at x). Then $\text{Lip}_p(f) = \sup_{x \in \mathbb{R}^n} \|J_f(x)\|_p$ where $\|J_f(x)\|_p$ is the induced operator norm on $J_f(x)$.

Therefore if f is a linear map represented by a matrix W when

$$\begin{aligned} \text{Lip}_p(f) &= \|W\|_p := \sup_{\|x\|_p=1} \|Wx\|_p \\ &= \begin{cases} \sigma_{\max}(W) & \text{if } p = 2 \\ \max_i \sum_j |W_{ij}| & \text{if } p = \infty \end{cases} \end{aligned}$$

where $\|W\|_p$ is the operator norm on matrices induced by the vector p -norm, and $\sigma_{\max}(W)$ is the largest singular value of W .

We can show the Lipschitz continuity of FC and CNN by applying the following lemma.

Lemma 1: ([17]) Let g, h be two composable Lipschitz functions. Then $g \circ h$ is also Lipschitz with $\text{Lip}(g \circ h) \leq \text{Lip}(g)\text{Lip}(h)$.

Corollary 1: For a fully-connected network or a convolutional neural network $f = W_k \circ \rho_{K-1} \circ W_{K-1} \circ \dots \circ \rho_1 \circ W_1$, we have $\text{Lip}_p(f) \leq \prod_k \|W_k\|_p$ under a choice of p -norm with 1-Lipschitz non-linearities ρ_k .

B. Attention Layers

The self-attention layer is not Lipschitz continuous [27]. However, when we consider the bounded inputs, the Lipschitz constant of attention layers $\text{Lip}_p(f)$ is bounded [27]. Since the input images are normalized and the network is trained with adamw [41] which penalizes the large weights, the inputs to attention layer is bounded. Therefore, we can assume the attention layer in our network is bounded with normalized image inputs.

C. Entire Transformer-based Network

Our transformer-based network is a composition of CNN, FC, and attention layers. By leveraging Lemma 1 on the normalized image space \mathcal{I} , we have

$$\text{Lip}_p(\text{Model}) \leq \text{Lip}_p(\text{CNN}) \cdot \text{Lip}_p(\text{FC}) \cdot \text{Lip}_p(\text{Attn})$$

where $\text{Lip}_p(\text{Model})$, $\text{Lip}_p(\text{CNN})$, $\text{Lip}_p(\text{FC})$ and $\text{Lip}_p(\text{Attn})$ are the Lipschitz constants of the entire model, all CNN layers, all FC layers and all attention layers respectively.

VII. AUDIO-VISUAL CORRESPONDENCE

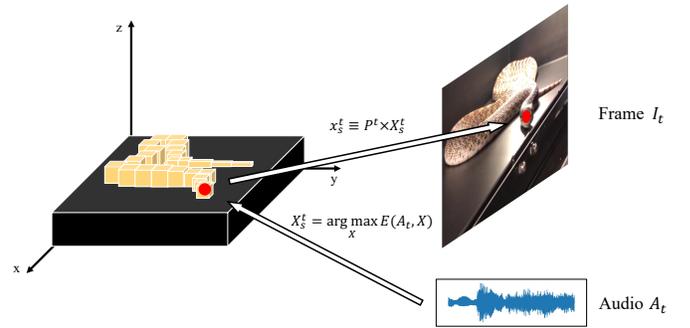


Fig. 11. Audio-visual relationship.

Given the audio recordings A^t , we can decode the 3D location of sound source by signal processing methods such as MUSIC algorithm [49] which searches for the location of the maximum of the spatial energy spectrum as

$$X_s^t = \arg \max_X E(A^t, X) \quad (9)$$

where $E(\cdot)$ is the spatial energy spectrum and $X_s^t \in \mathbb{R}^3$ is the 3D coordinate of sound source at time t . Therefore, the sound source can be easily corresponded to the pixel on a frame by applying camera matrix. Given that, we can represent the relationship between audio recordings A^t and 2D location of sound source x_s in frame I_t by

$$x_s^t \equiv P^t \times \arg \max_X E(A^t, P) \quad (10)$$

where P^t is the camera matrix at time t and x_t is the homogeneous coordinate of sound source in frame I_t . In practice, the camera matrix P^t is unknown but can be estimated by a adjacent input frames $P^t = \mathcal{H}(I_t, I_{t-1})$. Consequently, we can link the sound source location p_s^t in frames I_t from audio A_t .

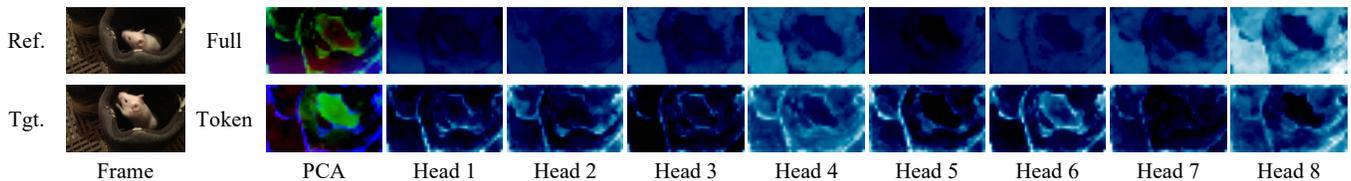


Fig. 12. **Visualization of reference-to-target attention map in last layer of transformer encoder.** We compare the attention map in the transformer encoder using full reference feature maps and compressed reference token. The PCA is conducted among the 8 heads of reference-to-target (r2t) attention map and keeps 3 main components for RGB channels. We only show the all 8 heads of the r2t attention. The attention map of token inputs is sharper on the instance region than that of full inputs which can also be verified by the PCA colormap.

VIII. DETAILED AUDIO PREPROCESSING

We first resample A_i to 16 kHz mono then compute the spectrum using magnitudes of the Short-Time Fourier Transform (STFT) with a window size of 25 ms, a window hop of 10 ms, and a periodic Hann window. We pad the audio sequence to ensure the output has the same length as the input. The mel spectrogram is computed by mapping the spectrogram to 64 mel bins covering the range 125-7500 Hz. A stabilized log mel spectrogram is further computed by applying $\log(\text{mel-spectrum} + 0.01)$ where the offset is used to avoid taking a logarithm of zero.

Method	Backbone	AP	AP50	AP75	AR1	AR10
CrossVIS	ResNet-50	42.2	59.1	47.4	42.1	49.9
CrossVIS	ResNet-101	44.7	60.4	50.5	42.0	49.7
Ours	ResNet-50	46.6	63.3	51.3	44.2	51.5
Ours	ResNet-101	48.6	63.9	52.5	44.6	52.9

TABLE XI
COMPARISON TO STATE-OF-THE-ART VIDEO INSTANCE SEGMENTATION ON AVIS VAL SET.

IX. ADDITIONAL EXPERIMENTS ON AVIS DATASET

We report the performance of using ResNet-50 and ResNet-101 backbone on AVIS dataset. As shown in Table XI, we also compare the CrossVIS baseline which only leverages the video data. Our method outperforms CrossVIS baseline by 4.4 and 3.9 mAP with ResNet-50 and ResNet-101 backbone respectively.

We also report the p-value of t-test on AVIS results with different thresholds. As shown in Table XII, the p-value under $\text{IoU} > 0.75$ threshold is smaller than that under $\text{IoU} > 0.5$ threshold, which means the audio inputs have more effect on the high-quality results.

X. DIFFERENCE WITH MASKFORMER [13]

MaskFormer [13] is an image-level panoptic segmentation method that treats prediction belonging to both stuff and things

Ref. number	AP	AP50	AP	AP50	p-value
	w/o Audio-token	w Audio-token	w Audio-token	w Audio-token	
1	42.9	60.2	44.8 (+1.9)	60.3	0.30 (0.24)
2	44.3	60.9	46.0 (+1.7)	63.8	0.26 (0.16)
3	45.6	61.9	46.6 (+1.0)	63.3	0.65 (0.42)
4	45.6	62.2	45.7 (+0.1)	62.7	-
5	45.2	62.9	45.0 (-0.2)	62.6	-

TABLE XII
IMPACT OF INTEGRATING AUDIO-TOKEN WITH DIFFERENT REFERENCE FRAMES ON AVIS VAL SET. THE P-VALUE IS COMPUTED BY T-TEST AMONG SAMPLES HAVING AN $\text{IoU} > 0.5$ ($\text{IoU} > 0.75$ IN THE BRACKET).

classes as a mask with corresponding semantic categories. The pixel and transformer decoder used in our model is similar to the structure adopted in MaskFormer. Here we formally analyze the difference between our method and MaskFormer.

Our method is a video-level instance segmentation method which equips with a context fusion module (CFM) to fuse the contextual information. The proposed CFM module is leveraged to compress the historical features from adjacent frames with low redundancy and construct the spatial correlation between target and contextual features using attention mechanism. The importance of the final segmentation is considered when fusing the context features. By utilizing the flexibility of CFM, we also fuse the audio modality to facilitate the instance segmentation task. The transformer decoder of our model takes the multi-modal features as inputs, which is different from MaskFormer. The decoded instance embedding is further to be leveraged to track the instance identities given the Lipschitz continuity of the network.

XI. VISUALIZATION OF ATTENTION MAPS IN CFM

We demonstrate the attention map of all eight heads of the reference-to-target attention. The attention map of the compressed feature is sharper than the attention map using the full features. The attention map is selected from the same position in the transformer encoder for fairly comparison.

XII. LIMITED GAIN OF INVOLVING AUDIO MODALITY

Our method introduces audio modality in an online fashion with synchronized the audio and video data. Since we considering streamlining video frames (as previous online VIS methods), we also introduce the corresponding audio frames in the same way. However, for audio data analysis, long-term historical information is required. Although we have reference audio frames, they are limited to a small number since we can only consider a moderate number of reference video frames with speed concerns. In this way, we believe that the audio modality can be further leveraged in asynchronous audio-visual inputs or clip-level inputs. We will investigate those settings in the future.