

# RETHINKING AUDIOVISUAL SEGMENTATION WITH SEMANTIC QUANTIZATION AND DECOMPOSITION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1        Audiovisual segmentation (AVS) is a challenging task that aims to segment visual  
 2        objects in videos based on their associated acoustic cues. With multiple sound  
 3        sources involved, establishing robust correspondences between audio and visual  
 4        contents poses unique challenges due to its (1) intricate entanglement across sound  
 5        sources and (2) frequent shift among sound events. Assuming sound events occur  
 6        independently, the multi-source semantic space (which encompasses all possible  
 7        semantic categories) can be represented as the Cartesian product of single-source  
 8        sub-spaces. This motivates us to decompose the multi-source audio semantics  
 9        into single-source semantics, enabling more effective interaction with visual con-  
 10        tent. Specifically, we propose a semantic decomposition method based on product  
 11        quantization, where the multi-source semantics can be decomposed and repre-  
 12        sented by several quantized single-source semantics. Furthermore, we introduce  
 13        a global-to-local quantization mechanism, which distills knowledge from stable  
 14        global (clip-level) features into local (frame-level) ones, to handle the constant  
 15        shift of audio semantics. Extensive experiments demonstrate that semantically  
 16        quantized and decomposed audio representation significantly improves AVS per-  
 17        formance, e.g., +21.2% mIoU on the most challenging AVS-Semantic benchmark.

## 18    1 INTRODUCTION

19        Recently, audiovisual segmentation (AVS) (Zhou et al., 2022) is introduced to explore audiovisual  
 20        correlations at the pixel level. Specifically, AVS aims to segment sounding object(s) in video frames  
 21        with the associated audio. Audiovisual semantic segmentation (AVSS) (Zhou et al., 2023) extends  
 22        AVS by additionally identifying the categories of sound sources. As shown in Fig. 1 (a), in contrast  
 23        to the visual domain, where each pixel has a unique semantic label, multi-source audio is tempo-  
 24        rally entangled, leading to ambiguity when associating the hybrid audio with semantically distinct  
 25        visual contents. This motivates us to explore suitable representations of multi-source audio for more  
 26        effective audiovisual interactions.

27        Let us commence by exploring the simplest scenario (left of Fig. 1 (a)), involving a single sound  
 28        source. In this scenario, each visual pixel or acoustic timestep is only associated with one semantic  
 29        label. Here, we denote the set containing all possible semantic labels for pixels as visual semantic

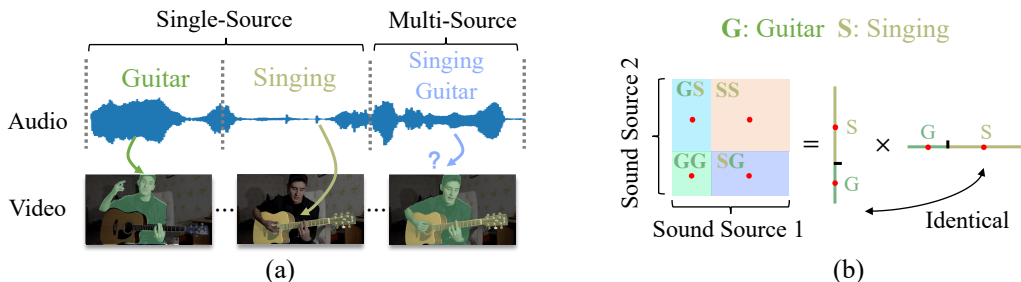


Figure 1: (a) **Audiovisual semantic interaction.** (b) **Semantic decomposition.** Multi-source audio semantic space can be assumed as a Cartesian product of single-source semantics, which can be decomposed via product quantization. The red points represent the quantized semantics.

space and those for acoustic timesteps as acoustic semantic space, respectively. In this example, we can find that both visual and acoustic single-source semantic spaces share the same semantic labels as {Guitar (G), Singing (S)}. Let us consider the two-source moment (right of Fig. 1 (a)). The visual semantics remain the same as in previous frames, but the size of possible two-source audio semantics presents a quadratic increase ({GG, GS, SG, SS}). This not only increases the difficulty in modeling larger semantic spaces but also complicates the alignment between visual and acoustic semantic spaces. The complexity further intensifies as more sources come into play.

Unlike previous methods (Zhou et al., 2023; 2022) that directly interact entangled multi-source audio representation with visual contents, we intend to disentangle the multi-source audio semantics into several single-source semantics for further more effective audiovisual interaction. We simplify the problem by assuming independent sound events, which allows us to represent the multi-source semantic space as a Cartesian product of identical single-source semantic spaces. In specific, we introduce a product quantization-based (PQ-based) method to decompose the multi-source semantics. Product quantization aims to represent a complex space through the product of several subspaces. In the multi-source case, single-source semantics can serve as subspaces. We show the semantic decomposition of a two-source example in Fig. 1 (b) where the single-source semantic subspaces share identical semantics {G, S}. Specifically, product quantization can be easily achieved by learning separate transforms of multi-source semantics and then quantizing them utilizing a shared codebook as shown in Fig. 2. We interact the decomposed single-source semantics with visual features for effective alignment.

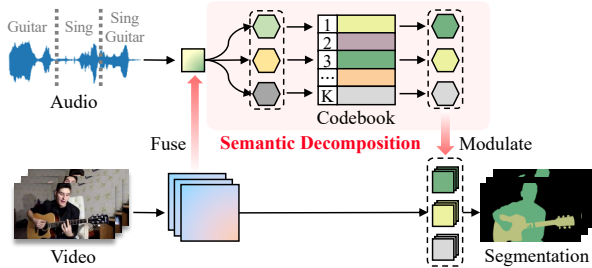


Figure 2: Semantic decomposition via product quantization (with sharing codebook for subspaces).

Furthermore, considering that active sound events may continually change over time, another challenge for AVS is to extract frame-level audio semantics which is typically not as robust as extracting from clip-level audio. To improve the frame-level audio representation, we propose a global-to-local mechanism, which distills knowledge from robust global (clip-level) audio representations into local (frame-level) ones. Specifically, we build an effective codebook for semantic quantization with clip-level visual-enriched audio features and then apply this codebook to perform local quantization on each frame without updating it. Thereby, the local semantic tokens are calibrated to the more robust and representative clip-level feature in the codebook.

In summary, our contribution is three-fold:

- An effective approach of multi-source audio semantic decomposition via product quantization, addressing the challenge of interacting visual and audio features in multiple object scenarios.
- A global-to-local distilling mechanism for frame-level audio semantic enhancement, addressing the ineffectiveness of frame-level audio feature extraction.
- Extensive experiments are conducted to verify the effectiveness of the proposed method, which significantly outperforms previous state-of-the-art methods on three AVS benchmarks, especially for multi-object datasets (+5.4% mIoU for AVS-Objec-Multi and +21.2% mIoU for AVS-Semantic).

## 2 RELATED WORK

**Audiovisual segmentation and localization.** Audiovisual segmentation (AVS), which was recently introduced (Zhou et al., 2022), aims to segment the objects that produce sound at the time of the image frame. Zhou et al. (Zhou et al., 2022) proposed a method with cross-modal attention to locate the sound source, making it the pioneering work in AVS. Recently, an extended task of AVS, audiovisual semantic segmentation (AVSS), is proposed by Zhou et al. (Zhou et al., 2023) which aims to not only segment the mask of sound sources but also predict the category of each sound

83 source. Due to the semantic entanglement in audio, tackling multi-source AVSS is more challenging  
 84 than AVS task. Zhou *et al.* (Zhou et al., 2023) follows the TPAVI module in (Zhou et al., 2022)  
 85 to conduct audiovisual interaction. Sound source localization (SSL) (Mo & Morgado, 2022a;b;  
 86 Senocak et al., 2018; Hu et al., 2019; Qian et al., 2020; Chen et al., 2021; Afouras et al., 2020) is  
 87 a related problem to AVS that aims to locate the regions of sounds in the visual frame. Common  
 88 SSL methods (Arandjelovic & Zisserman, 2018; 2017; Cheng et al., 2020; Senocak et al., 2018)  
 89 leverage cross-modal correspondence between audio and visual features to locate sounds, which are  
 90 then displayed as heatmaps. For instance, Mo *et al.* (Mo & Morgado, 2022a) leverage multi-level  
 91 audiovisual contrastive learning to effectively locate the objects. Different from previous methods  
 92 primarily designed for single-source scenarios, our objective is to address the semantic entanglement  
 93 present in multi-source audios and explore methods for effective interaction between multi-source  
 94 audios and videos.

95 **Audio-visual learning.** Audio-visual learning has been explored in many works (Aytar et al., 2016;  
 96 Arandjelovic & Zisserman, 2017; Korbar et al., 2018; Senocak et al., 2018; Zhao et al., 2018; 2019;  
 97 Gan et al., 2020; Georgescu et al., 2022) which aims to learn audio-visual correspondence from  
 98 paired audio-visual data. Most methods maximize the mutual information between corresponding  
 99 audio and video pairs by several proxy tasks. Constructing negative samples (Zhao et al., 2018;  
 100 2019; Gan et al., 2020) and learning to push them away while closing positive ones is a common  
 101 goal. Recently, another track (Georgescu et al., 2022; Gong et al., 2022) masks information in  
 102 audio-visual pairs and tries to reconstruct incomplete information in one modality by conditioning  
 103 on the other. The learned correspondence can be leveraged for several tasks, such as audio-visual  
 104 source localization (Mo & Morgado, 2022a;b; Senocak et al., 2018; Hu et al., 2019; Qian et al.,  
 105 2020; Chen et al., 2021; Afouras et al., 2020), audio-visual separation (Gao & Grauman, 2019;  
 106 Morgado et al., 2018; 2020; Chen et al., 2020a), audio-visual parsing (Wu & Yang, 2021; Mo &  
 107 Tian, 2022; Lin et al., 2021; Tian et al., 2020). In this work, we focus on how to effectively construct  
 108 correspondence between multi-source audio and video for fine-grained audiovisual segmentation  
 109 which is more challenging due to the entanglement of semantics in audio.

### 110 3 METHOD

111 In this section, we first present the formulation of the product quantization-based (PQ-based) method  
 112 for multi-source audio semantic decomposition. Then, we outline the pipeline that utilizes the quan-  
 113 tized and decomposed audio representation to improve the audiovisual segmentation tasks.

#### 114 3.1 PQ-BASED MULTI-SOURCE SEMANTIC DECOMPOSITION.

115 The core of the PQ-based decomposition is to concisely represent the multi-source semantic space  
 116  $\mathcal{X}_m$  with the product of multiple single-source semantic spaces  $\mathcal{X}_s$ .

117 Given a codebook containing a finite set of codewords  $\mathcal{C} = \{e_i\}_{i=1}^K$ , the vector quantizer  $VQ(\cdot)$   
 118 maps a feature  $x \in \mathcal{X}$  to a codeword  $e_i = \arg \min_{e_i \in \mathcal{C}} \|x - e_i\|_p$  that minimizes the distance  
 119 between  $x$  and  $e_i$  in the  $p$ -norm sense. As the single-source audio is semantically unique for each  
 120 time step, for single-source space with  $K$  sound event categories, a codebook  $\mathcal{C}_s$  with  $K_s = K$   
 121 codewords can sufficiently encode the space  $\mathcal{X}_s$  without losing information. Nevertheless, for a  $N$ -  
 122 source semantic space, a combination of sound events can appear for each time step. Therein, to  
 123 fully represent the space, a codebook  $\mathcal{C}_m$  of size  $K_m = K^N$  is required.

124 We assume a  $N$ -source semantic space  $\mathcal{X}_m$  is a Cartesian product of several identical single-source  
 125 semantic spaces  $\mathcal{X}_s$  as

$$\mathcal{X}_m = \underbrace{\mathcal{X}_s \times \cdots \times \mathcal{X}_s}_N. \quad (1)$$

126 Specifically, we can obtain the product quantization of  $x \in \mathcal{X}_m$  through an order-invariant operation,  
 127 *i.e.*, concatenation, on separately quantized  $x_i = f_i(x)$ , where  $f_i(\cdot)$  is a transform applied on  $x$ :

$$PQ(x) = VQ_s(x_1) \oplus \cdots \oplus VQ_s(x_N), \quad \text{w.r.t } VQ_s \sim \mathcal{C}_s. \quad (2)$$

128  $VQ_s \sim \mathcal{C}_s$  denotes the  $VQ_s(\cdot)$  is associated with codebook  $\mathcal{C}_s$  and  $\oplus$  denotes channel-wise concate-  
 129 nation. As the codebook  $\mathcal{C}_s$  is shared with all  $VQ_s(\cdot)$ , the codebook for the multi-source semantic

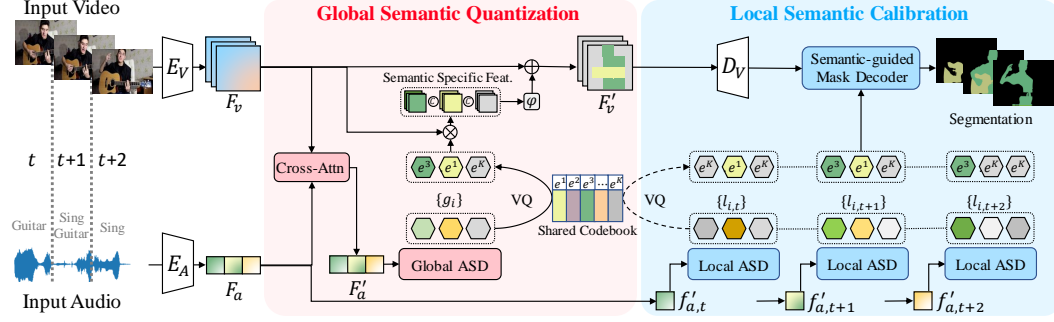


Figure 3: Method overview. The **global semantic quantization** module decomposes multi-source audio features and enables interaction between the decomposed single-source audio feature and visual features. The **local semantic calibration** module distills knowledge from global (clip-level) audio features to local (frame-level) audio features by utilizing a shared codebook, which stores quantized audio representation during semantic quantization.

130 space  $\mathcal{C}_m$  is reduced to be the same as  $\mathcal{C}_s$ . By constraining the size of the single-source codebook  
 131  $K_s \ll K^N$ , we can force the transform  $f_i(\cdot)$  to decompose the multi-source semantics  $x$ .

### 132 3.2 NETWORK OVERVIEW

133 We present the proposed framework with **Semantically Quantized and Decomposed (SQD)** audio  
 134 representation consisting of three main components: feature encoding, global semantic quantization  
 135 and local semantic calibration, as illustrated in Fig. 3.

136 (1) First, we extract visual features  $F_v = \{f_{v,t}\}_{t=1}^T$  and acoustic features  $F_a = \{f_{a,t}\}_{t=1}^T$  by sepa-  
 137 rate encoders. (2) Then, **to decompose semantics in multi-source audio features**, we use a global  
 138 semantic decomposition module to map the audio query into a set of semantic tokens  $\{g_i\}_{i=1}^N$ . We  
 139 then learn a semantic codebook to quantize them. The quantized tokens are further employed to  
 140 modulate the visual features to inject information about corresponding sound sources. (3) After-  
 141 wards, **to obtain frame-level audio features to query object masks**, we utilize a local semantic  
 142 decomposition module for each time step, which uses the global codebook to decouple local audio  
 143 semantics. Each quantized local semantic token  $VQ(l_{i,t})$  serves as a query to segment a frame-level  
 144 mask with the semantic-guided mask decoder. Overall, the proposed SQD boosts AVS by enhancing  
 145 the audiovisual semantic interaction.

### 146 3.3 GLOBAL SEMANTIC QUANTIZATION

147 To tackle the mixture of multi-source audio queries and effectively conduct audiovisual fusion, we  
 148 propose global semantic quantization to decompose audio semantics, which consists of two steps:  
 149 global semantic decomposition and audiovisual semantic recombination. The detailed structure of  
 150 the modules is illustrated in Fig. 4.

151 **Global semantic decomposition.** Global semantic decomposition aims to decompose multi-source  
 152 audio semantics into single-source semantics. Specifically, the audio feature  $F_a$  is first fused with  
 153 video feature  $F_v$  to be  $F'_a$ , taking the form:

$$\begin{aligned} F'_a &= \text{LN}(\text{FFN}(h_a) + h_a), \\ h_a &= \text{LN}(\text{MCA}(F_a, F_v) + F_a), \end{aligned} \quad (3)$$

154 where MCA denotes Multi-head Cross-Attention, LN denotes Layer Normalization, and FFN de-  
 155 notes Feed-Forward Network. After that, we transform the audio feature  $F'_a$  to  $N$  decomposed  
 156 semantic tokens  $\{g_i\}_{i=1}^N$  with a global audio semantic decoder (global ASD):

$$g_i = \text{TrD}_{\text{global}}(p_i | F'_a) \quad (4)$$

157 by querying a set of learnable semantic prototypes  $\{p_i\}_{i=1}^N$  to the feature  $F'_a$ , with a transformer  
 158 decoder  $\text{TrD}_{\text{global}}$ . Each semantic token is then quantized to be  $e_i = VQ(g_i)$  with the shared  
 159 codebook  $\mathcal{C} = \{e^k\}_{k=1}^K$ , imposing that all semantic tokens to share an identical feature subspace

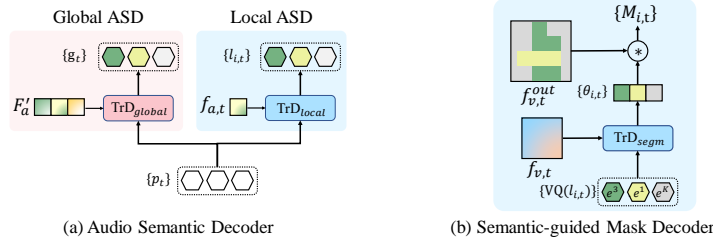


Figure 4: (a) Global and local audio semantic decoder (ASD) share similar structures that query clip-/frame-level audio features,  $F'_a$  or  $f_{a,t}$ , with a transformer decoder  $\text{TrD}_{global}/\text{TrD}_{local}$  using learnable semantic prototypes  $\{p_i\}$ . (b) The semantic-guided mask decoder contains a transformer decoder  $\text{TrD}_{segm}$  to align audiovisual features and computes dynamic filters  $\theta_{i,t}$ . The final mask  $M_{i,t}$  is generated by a dynamic convolution between the visual feature  $f_{v,t}^{out}$  and  $\theta_{i,t}$ .

160 with low cardinality. Note that we set the codebook size  $K \ll D_{semantic}^N$  to force the network to  
 161 learn decomposed semantics, where  $D_{semantic}^N$  is the number of sound event categories  $D_{semantic}$   
 162 to the power of  $N$ .

163 **Audiovisual semantic recombination.** Audiovisual semantic recombination aims to leverage the  
 164 decomposed audio feature to interact with visual features. After obtaining quantized global semantic  
 165 tokens  $\{e_i\}_{i=1}^N$ , which encode  $N$  groups of decomposed semantics, we aim to interact them with  
 166 visual features while preserving the original function of the multi-source audio input. A set of  
 167 dynamic filters  $\{w_i \in \mathbb{R}^{C_v}\}_{i=1}^N$  are first learned from global semantic tokens  $\{g_i\}_{i=1}^N$  by two linear  
 168 layers. After that, we utilize channel-wise attention to modulate video features by each filter to  
 169 interact the visual feature with the content referred by different semantic tokens, which is given by:

$$F'_v = \text{BN}(\varphi(w_i F_v \oplus \dots \oplus w_N F_v) + F_v), \quad (5)$$

170 where  $\varphi$  denotes a convolution layer to reduce channel from  $N \times C_v$  to  $C_v$ , BN denotes Batch Nor-  
 171 malization, and  $\oplus$  denotes concatenation among channels. By incorporating channel-wise attention,  
 172 the visual features can be more effectively concentrated on the relevant audio content. Furthermore,  
 173 through channel-wise concatenation, the decomposed audio semantics can be reintegrated, produc-  
 174 ing hybrid semantics that refers to the holistic contents of the original audio input.

### 175 3.4 LOCAL SEMANTIC CALIBRATION

176 Since the audio query is time-variant, global semantic tokens cannot be accurately aligned with  
 177 visual features at the frame level. To segment audio-queried contents in each frame, we propose the  
 178 local semantic calibration, consisting of a local semantic decomposition stage and a semantic-guided  
 179 mask decoding stage.

180 **Local semantic decomposition.** Local semantic decomposition module aims to decompose the  
 181 semantics encoded in each audio frame. Similar to the global semantic decoder, the local semantic  
 182 decoder (Local ASD) decodes frame-level semantics with a transformer decoder  $\text{TrD}_{local}$  and a set  
 183 of semantic prototypes  $\{p_i\}_{i=1}^N$ . The local semantic tokens  $l_{i,t}$  are given by

$$l_{i,t} = \text{TrD}_{local}(p_i | f_{a,t}). \quad (6)$$

184 The local semantic tokens do not build their own codebook but utilize the global codebook  $\mathcal{C}$ , that  
 185 is, they do not update  $\mathcal{C}$  but are committed to being close to the vectors in  $\mathcal{C}$ . In this way, the local  
 186 semantic tokens distill knowledge from the global ones. Further explanation regarding supervision  
 187 will be provided in Section 3.5.

188 **Semantic-guided mask decoding.** We utilize the semantic-guided mask decoder to decode visual  
 189 features into masks that correspond to decomposed local audio semantics, with detailed structure  
 190 illustrated in Fig. 4 (b). Pyramid video features  $F_v^{out} = \{f_{v,t}^{out}\}_{t=1}^T$  are obtained with the feature  
 191 pyramid network (Lin et al., 2017a). We leverage a shared multimodal transformer decoder  $\text{TrD}_{segm}$   
 192 to generate dynamic filters  $\theta_{i,t} = \phi_{segm}(\text{TrD}_{segm}(\text{VQ}(l_{i,t}) | f_t))$  for each timestep, where  $\phi_{segm}$  is  
 193 a two-layer fully-connected network. The final mask segmentation can be obtained by:

$$M_{i,t} = f_t^{out} * \theta_{i,t}, \quad (7)$$



Method	Backbone	AVS-Object-Single			AVS-Object-Multi			AVS-Semantic
		$\mathcal{J}\&\mathcal{F}\uparrow$	$\mathcal{J}\uparrow$	$\mathcal{F}\uparrow$	$\mathcal{J}\&\mathcal{F}\uparrow$	$\mathcal{J}\uparrow$	$\mathcal{F}\uparrow$	mIoU $\uparrow$
ResNet Backbone								
LVS (Chen et al., 2021)	ResNet-18	44.5	37.9	51.9	31.3	29.5	33.0	-
MSSL (Qian et al., 2020)	ResNet-18	55.6	44.9	66.3	31.4	26.1	36.3	-
3DC (Mahadevan et al., 2020)	3DC	66.5	57.1	75.9	43.6	36.9	50.3	17.3
AOT (Yang et al., 2021)	ResNet-50	-	-	-	-	-	-	25.4
AVS (Zhou et al., 2023)	ResNet-50	78.8	72.8	84.8	53.6	47.9	57.8	20.2
Bi-Gen (Hao et al., 2023)	ResNet-50	79.8	74.1	85.4	50.9	50.0	56.8	-
AVSegFormer (Gao et al., 2023)	ResNet-50	81.2	76.5	85.9	56.2	49.5	62.8	24.9
<b>SQD (Ours)</b>	ResNet-50	<b>81.8</b>	<b>77.6</b>	<b>86.0</b>	<b>61.6</b>	<b>59.6</b>	<b>63.5</b>	<b>46.6</b>
Transformer Backbone								
iGAN (Mao et al., 2021)	Swin-Base*	69.7	61.6	77.8	48.7	42.9	54.4	-
SST (Duke et al., 2021)	SSL	73.2	66.3	80.1	49.9	42.6	57.2	-
LGVT (Zhang et al., 2021)	Swin-Base*	81.1	74.9	87.3	50.0	40.7	59.3	-
AVS (Zhou et al., 2023)	PVT-v2-Base	83.3	78.7	87.9	59.3	54.0	64.5	29.8
<b>SQD (Ours)</b>	Swin-Tiny	<b>83.9</b>	<b>79.5</b>	<b>88.2</b>	<b>64.0</b>	<b>61.9</b>	<b>66.1</b>	<b>53.4</b>
<b>SQD (Ours)</b>	V-Swin-Tiny	<b>84.7</b>	<b>80.7</b>	<b>88.7</b>	<b>65.4</b>	<b>63.7</b>	<b>67.0</b>	<b>54.7</b>

Table 1: **Quantitative comparison to AVS and AVSS methods.** Swin-Base\* denotes modified Swin-Base Transformer (Liu et al., 2021). SSL is Sparse Spatiotemporal Transformers (Duke et al., 2021). PVT-v2 (Wang et al., 2022) is a strong Pyramid Vision Transformer. V-Swin-Tiny is the Video Swin Transformer (Liu et al., 2022).  $\uparrow$  indicates the larger the better.

194 where \* denotes the dynamic convolution (Chen et al., 2020b). Each filter represents semantics of  
 195 a decomposed single-source audio, contributing to the segmentation of the single sounding object.  
 196 Additional class probability prediction  $P_{i,t}$  and bounding box prediction  $B_{i,t}$  for each mask  $M_{i,t}$  are  
 197 performed by two two-layer fully connected networks from the output of  $\text{TrD}_{\text{segm}}(\text{VQ}(l_{i,t})|f_t)$ .

### 198 3.5 LOSS FUNCTION

199 The overall loss function is given by

$$\mathcal{L} = \lambda_{\text{quant}}\mathcal{L}_{\text{quant}} + \mathcal{L}_{\text{segm}}, \quad (8)$$

200 where  $\mathcal{L}_{\text{quant}}$  and  $\mathcal{L}_{\text{segm}}$  are the loss for semantic quantization and semantic segmentation, respec-  
 201 tively.  $\lambda_{\text{quant}}$  is a constant.

202 **Loss for semantic quantization.** The quantizer is shared with both global and local semantic de-  
 203 composition, while the local semantic tokens do not update the codebook. The loss is given by

$$\begin{aligned} \mathcal{L}_{\text{quant}} = & \sum_{i=1}^N \|\text{VQ}(g_i) - \text{sg}[g_i]\|_2^2 \\ & + \lambda_{\text{com}} \|\text{sg}[\text{VQ}(g_i)] - g_i\|_2^2 + \lambda_{\text{com}} \|\text{sg}[\text{VQ}(l_i)] - l_i\|_2^2, \end{aligned} \quad (9)$$

204 where  $\text{sg}[\cdot]$  stands for stop-gradient operation.  $\text{VQ}(\cdot)$  denotes the vector quantization function,  
 205 where  $\text{VQ}(x) = e_i = \arg \min_{e_i} \|x - e_i\|_2 \in \mathcal{C}$  and  $\mathcal{C} = \{e_i\}_{i=1}^K$  is the shared codebook. The first  
 206 term aims to update the codebook. The second and third terms aim to minimize the quantization  
 207 error by forcing the input vector to be quantized to its closest vector in the codebook.

208 **Loss for semantic segmentation.** Let the predictions of the network be  $\mathbf{y} = \{y_i\}_{i=1}^N$  where  
 209  $y_i = \{B_{i,t}, P_{i,t}, M_{i,t}\}_{t=1}^T$ .  $B_{i,t}$ ,  $P_{i,t}$  and  $M_{i,t}$  denote bounding box, class probability and mask  
 210 predictions respectively. We denote the ground-truth as  $\hat{\mathbf{y}} = \{\hat{y}_j\}_{j=1}^N$  (padded with  $\emptyset$  (Cheng et al.,  
 211 2021a)) where  $\hat{y}_j = \{\hat{B}_{j,t}, \hat{C}_{j,t}, \hat{M}_{j,t}\}_{t=1}^T$ .  $C_{j,t}$  is the ground-truth class for the  $j$ -th sounding object  
 212 in the video at  $t$  frame. We search for an assignment  $\sigma \in \mathcal{P}_N$  with the highest similarity where  $\mathcal{P}_N$   
 213 is a set of permutations of  $N$  elements. The similarity can be computed as

$$\mathcal{L}_{\text{match}}(y_i, \hat{y}_j) = \lambda_{\text{box}}\mathcal{L}_{\text{box}} + \lambda_{\text{cls}}\mathcal{L}_{\text{cls}} + \lambda_{\text{mask}}\mathcal{L}_{\text{mask}}, \quad (10)$$

214 where  $\lambda_{\text{box}}$ ,  $\lambda_{\text{cls}}$ , and  $\lambda_{\text{mask}}$  are weights to balance losses. We leverage a combination of Dice (Li  
 215 et al., 2019) and BCE loss as  $\mathcal{L}_{\text{mask}}$ , focal loss (Lin et al., 2017b) as  $\mathcal{L}_{\text{cls}}$ , and GloU (Rezatofighi  
 216 et al., 2019) and L1 loss as  $\mathcal{L}_{\text{box}}$ . The best assignment  $\hat{\sigma}$  is solved by the Hungarian algorithm (Kuhn,  
 217 1955). Given the best assignment  $\hat{\sigma}$ , the segmentation loss between ground-truth and predictions is  
 218 defined as  $\mathcal{L}_{\text{segm}} = \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\hat{\sigma}(j)})$ .



Figure 5: Qualitative comparison to Zhou et al. (Zhou et al., 2023) on AVS-Semantic. Each color represents a semantic category. Note that the class labels in the first row serve as references but are not given in the input.

219 4 EXPERIMENTS

220 **Dataset.** We conduct experiments on AVS-Object (Zhou et al., 2022) for AVS task and AVS-  
 221 Semantic (Zhou et al., 2023) for AVSS task.

- 222 • **AVS-Object:** AVS-Object dataset contains 5,356 short videos with corresponding audios in  
 223 which 4,932 audios contain single-source and 424 audios contain multiple sources. Class-  
 224 agnostic masks are given as annotations for AVS task. Typically, it is evaluated separately  
 225 for single- and multi-source audios as AVS-Object-Single and AVS-Object-Multi.
- 226 • **AVS-Semantic:** AVS-Semantic is an extended dataset from AVS-Object which contains  
 227 12,356 videos with 70 classes. Semantic segmentation is annotated for AVSS task. Both  
 228 single- and multi-source audio cases exist in the AVS-Semantic.

229 **Metrics.** For AVS task, the convention is to compute region similarity  $\mathcal{J}$  and contour accuracy  $\mathcal{F}$   
 230 as defined in (Pont-Tuset et al., 2017). Note that we follow the video segmentation convention to  
 231 use the region similarity  $\mathcal{J}$ , which is equivalent to mIoU in the binary AVS setting. For AVSS, we  
 232 follow the semantic segmentation convention to evaluate the model using mIoU which is defined as  
 233 the intersection over union averaged among all classes.

234 **Implementation Detail.** We implement our method in PyTorch (Paszke et al., 2019). We train our  
 235 model for 13 epochs and 16 epochs with a learning rate multiplier of 0.1 at the 11<sup>th</sup> and 14<sup>th</sup> epochs  
 236 for AVS-Object and AVS-Semantic datasets, respectively. The initial learning rate is 1e-4, and a  
 237 learning rate multiplier of 0.5 is applied to the backbone. We adopt batchsize = 4 and an AdamW  
 238 (Loshchilov & Hutter, 2017) optimizer with weight decay  $5 \times 10^{-4}$ . Multi-scale training is adopted  
 239 to obtain a strong baseline, and if no specification, all images are resized to have the longest side  
 240 224 during evaluation. More details are available in the supplementary materials.

241 4.1 MAIN RESULTS

242 **Quantitative comparison on AVS-Object.** Our method outperforms the previous state-of-the-art  
 243 (SOTA) method AVSegFormer (Gao et al., 2023) by 0.6 and 5.4 of  $\mathcal{J}$  &  $\mathcal{F}$  score on AVS-Object-  
 244 Single and AVS-Object-Multi datasets respectively (with ResNet-50 backbone). We notice that  
 245 the improvement on the multi-source setting is much larger than the single-source setting. This is  
 246 because single-source audios contain simple and disentangled semantics and can be easily aligned  
 247 with visual features while, for multi-source audios, the complex semantic space makes the alignment  
 248 to visual contents much more difficult.

249 **Quantitative comparison on AVS-Semantic.** Compared to the AVS-Object task, our method  
 250 demonstrates greater improvement in the AVS-Semantic task. As shown in the Table 1, our method  
 251 eclipses the previous SOTA AVSS method AOT (Yang et al., 2021) by a remarkable 21.2 mIoU  
 252 with ResNet-50 backbone. The improvement in the AVSS task can be attributed to several factors.  
 253 First, the task itself involves the semantic prediction of sound sources. However, due to mixed audio

Module	AVS-Object-Multi			AVS-Semantic
	$\mathcal{J}\&\mathcal{F}\uparrow$	$\mathcal{J}\uparrow$	$\mathcal{F}\uparrow$	mIoU $\uparrow$
Baseline	52.9	50.1	55.7	33.5
+GSD	56.7	54.5	58.8	38.4
+GSD+AVSR	58.6	56.5	60.6	40.9
+GSD+AVSR+LSD	60.1	58.2	61.9	44.5
+GSD+AVSR+LSD+SC	61.6	59.6	63.5	46.6

Table 2: Component analysis. GSD: global semantic decomposition; AVSR: audiovisual semantic recombination; LSD: local semantic decomposition; SC: sharing codebook.

Codebook Size	Object-M $\mathcal{J}\&\mathcal{F}\uparrow$	Semantic mIoU $\uparrow$
1	52.7	24.8
32	61.4	31.5
64	60.0	43.2
128	61.6	46.6
256	60.6	46.1

Table 3: Ablation on codebook size.

Token Number	Object-M $\mathcal{J}\&\mathcal{F}\uparrow$	Semantic mIoU $\uparrow$
1	59.7	40.2
3	61.0	43.5
5	61.6	46.6
7	61.6	45.9
9	61.2	46.3

Table 4: Ablation on decomposed token number.

Decomp. Domain	Object-M $\mathcal{J}\&\mathcal{F}\uparrow$	Semantic mIoU $\uparrow$
Time	56.2	38.9
Semantic	61.6	46.6

Table 5: Ablation on sound decomposition domain.

signals, aligning visual content accurately becomes challenging, leading to difficulties in classification. Secondly, the number of sound sources and categories of AVS-Semantic dataset are larger than AVS-Object, which will result in a larger complex semantic space. When the mixed semantics are not decomposed, the network struggles to handle the numerous mixed semantics effectively. Thirdly, in the AVS-Semantic dataset, sound event changes occur more frequently. As a result, a more robust frame-level audiovisual correspondence is required. Our proposed global-to-local distilling mechanism addresses this challenge by enhancing the capture of local semantic information, enabling accurate object segmentation.

**Qualitative comparison.** As shown in Fig. 5, we qualitatively compare our method to the method proposed by Zhou et al. (Zhou et al., 2022) on AVS-Semantic. Our method achieves better results on both segmenting quality and class prediction accuracy. Since the method (Zhou et al., 2022) directly fuses mixed audio features with video features, we notice that it suffers from object incorrectness when multiple sound sources are present. Meanwhile, due to the lack of frame-level audiovisual calibration, (Zhou et al., 2022) cannot effectively handle the audio semantic changes. More qualitative results on the AVS-Object dataset are available in the Appendix.

## 4.2 ABLATION STUDY

**Module Effectiveness.** We conduct experiments to validate the effectiveness of our proposed modules. We first construct a **baseline** with unimodal encoders and the semantic-guided mask decoder, and then add other modules step-by-step. As shown in Table 2, each of the proposed modules benefits the performance. For AVS-Semantic, both global semantic decomposition (GSD) and local semantic decomposition (LSD) bring obvious gains; for AVS-object, the LSD only slightly improves the performance. This could be attributed to the longer duration of videos in AVS-Semantic compared to AVS-Object, which allows for a greater number of semantic changes within each clip. Finally, with all components, our method achieves the gains of 10.5  $\mathcal{J}\&\mathcal{F}$  and 13.1 mIoU on AVS-Object-Multi and AVS-Semantic respectively when compared to the baseline.

**Semantic token number.** We ablate the semantic token number for the global-ASD and local-ASD in Table 4. We observe that a token number of 5 yielded the best performance. This can be attributed to the fact that the maximum number of mixed sound sources for the audio-visual dataset is 5.

**Codebook size.** The cardinality of the codebook is essential to our semantic decomposition. Ideally, we aim to constrain the cardinality of the codebook to be close to the semantic category number. We ablate on codebook size from 1 to 256. When the codebook size equals 1, all the decomposed audio tokens are the same, resulting in all the same segmentation results. As shown in Table 3, we notice even a codebook size of 1 achieves 24.8 mIoU on AVS-Semantic. A codebook of size=128 achieves the best performance. Please note that a codebook size slightly larger than the category number, e.g.128, will not hamper the semantic decomposition capability of our method, since  $128 \ll 70^N$  where  $N > 1$  is the maximum sound source number, and 70 is the category number.



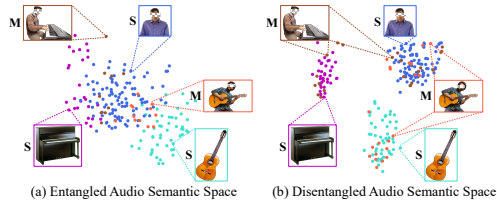


Figure 6: Visualization comparison between entangled and our disentangled audio semantic space. “M” and “S” notations denote multi-source and single-source inputs.

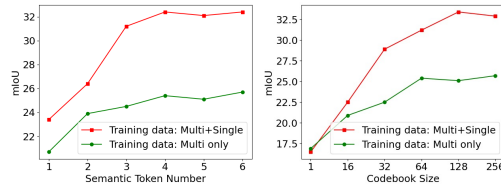


Figure 7: Comparison of training w. and w/o. single-source data.

### 290 4.3 ANALYSIS

291 **Visualization of decomposed semantic space.** As shown in Fig. 6, we visualize the semantic  
 292 space with and without semantic decomposition on the AVS-Semantic dataset using t-SNE (Van der  
 293 Maaten & Hinton, 2008). Three types of single-source audios (“man”, “guitar”, “piano”) and two  
 294 types of multi-source audios (“man+guitar”, “man+piano”) are enrolled. Without decomposition,  
 295 the multi-source features are highly entangled, presenting fewer evidences related to single-source  
 296 semantics. However, after performing semantic decomposition, the “man+guitar” feature presents  
 297 clear evidences related to its corresponding single-source (“man” and “guitar”) semantics. This is  
 298 reflected in the proximity of the ”man+guitar” feature to the centroids of its corresponding single-  
 299 source features. The same applies to the ”man+piano” feature. Note that, we omit the “background”  
 300 feature in the visualization.

301 **Importance of the single-source audio on the semantic decomposition of multi-source audio**  
 302 **representation.** We present empirical evidence that the single-source audio samples significantly  
 303 contribute to the success of semantic decomposition. To demonstrate this, we compare the perform-  
 304 ance of our model trained on two training sets with the same number of samples: one contains  
 305 solely multi-source audio samples, and the other contains single- and multi-source audio samples  
 306 with the ratio of 1:1. As illustrated in Fig. 7, the model trained solely on multi-source audio samples  
 307 exhibits inferior performance compared to the model trained on both types of samples, regardless of  
 308 the token number and codebook size. We conjecture that the single-source samples serve as infor-  
 309 mative anchors that assist the model in learning the correct distributions of the decomposed simplex  
 310 spaces for multi-source samples. In the absence of single-source samples, the decomposition task  
 311 could be more difficult due to the absence of such informative anchors.

312 **Ablation on audio decomposition domain.** We conducted an experiment to demonstrate the bene-  
 313 fits of conducting audio decomposition at the semantic domain instead of the time domain. Specifi-  
 314 cally, we decomposed the multi-source audio with a commonly used sound source separation model  
 315 (Chen et al., 2022) and then performed audiovisual segmentation for each decomposed audio using  
 316 our proposed model. The results in Table 5 clearly show that our semantic-level decomposition  
 317 mechanism outperforms the time-domain decomposition approach. We attribute this improvement  
 318 to two factors: 1) the imperfection of sound source separation and 2) the conflicts that arise when  
 319 combining the masks for each source in the time domain without considering visual content during  
 320 separation. In contrast, our semantic-domain approach does not suffer from these issues and can  
 321 effectively leverage the information contained in both audio and visual modalities.

## 322 5 CONCLUSION

323 This paper presents an approach to address the challenges in audiovisual segmentation by proposing  
 324 semantic decomposition of complex semantic spaces that encode multi-source audios, followed by  
 325 their interaction with visual features. This reduces the semantic ambiguity in multi-source audio-  
 326 visual interaction. To handle sound event changes, we propose local semantic calibration to align  
 327 audio and video on a per-frame basis. Our method also incorporates a codebook sharing mechanism  
 328 to enhance local audio features by distilling knowledge from that at the global level. The proposed  
 329 approach is evaluated on three AVS benchmarks and the results demonstrate its superiority and ef-  
 330 fectiveness over previous state-of-the-art methods.

## 331 REFERENCES

- 332 Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised  
333 learning of audio-visual objects from video. In *Computer Vision–ECCV 2020: 16th Euro-*  
334 *pean Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pp. 208–224.  
335 Springer, 2020. 3
- 336 Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE*  
337 *International Conference on Computer Vision*, pp. 609–617, 2017. 3
- 338 Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *Proceedings of the European*  
339 *conference on computer vision (ECCV)*, pp. 435–451, 2018. 3
- 340 Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from  
341 unlabeled video. *Advances in neural information processing systems*, 29, 2016. 3
- 342 Adam Botach, Evgenii Zheltonozhskii, and Chaim Baskin. End-to-end referring video object seg-  
343 mentation with multimodal transformers. In *Proceedings of the IEEE/CVF Conference on Com-*  
344 *puter Vision and Pattern Recognition*, pp. 4985–4995, 2022. 15
- 345 Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah,  
346 Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. Soundspaces: Audio-visual nav-  
347 igation in 3d environments. In *Computer Vision–ECCV 2020: 16th European Conference, Glas-*  
348 *gow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pp. 17–36. Springer, 2020a. 3
- 349 Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zis-  
350 serman. Localizing visual sounds the hard way. In *Proceedings of the IEEE/CVF Conference on*  
351 *Computer Vision and Pattern Recognition*, pp. 16867–16876, 2021. 3, 6
- 352 Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Zero-  
353 shot audio source separation through query-based learning from weakly-labeled data. In *Pro-*  
354 *ceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 4441–4449, 2022.  
355 9
- 356 Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic  
357 convolution: Attention over convolution kernels. In *Proceedings of the IEEE/CVF Conference on*  
358 *Computer Vision and Pattern Recognition*, pp. 11030–11039, 2020b. 6
- 359 Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexan-  
360 der G Schwing. Mask2former for video instance segmentation. *arXiv preprint arXiv:2112.10764*,  
361 2021a. 6
- 362 Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Modular interactive video object segmentation:  
363 Interaction-to-mask, propagation and difference-aware fusion. *arXiv preprint arXiv:2103.07941*,  
364 2021b. 15
- 365 Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved  
366 memory coverage for efficient video object segmentation. *arXiv preprint arXiv:2106.05210*,  
367 2021c. 15
- 368 Ying Cheng, Ruize Wang, Zhihao Pan, Rui Feng, and Yuejie Zhang. Look, listen, and attend: Co-  
369 attention network for self-supervised audio-visual representation learning. In *Proceedings of the*  
370 *28th ACM International Conference on Multimedia*, pp. 3884–3892, 2020. 3
- 371 Brendan Duke, Abdalla Ahmed, Christian Wolf, Parham Aarabi, and Graham W Taylor.  
372 Sstvos: Sparse spatiotemporal transformers for video object segmentation. *arXiv preprint*  
373 *arXiv:2101.08833*, 2021. 6
- 374 Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. Music gesture  
375 for visual sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
376 *and Pattern Recognition*, pp. 10478–10487, 2020. 3
- 377 Ruohan Gao and Kristen Grauman. 2.5 d visual sound. In *Proceedings of the IEEE/CVF Conference*  
378 *on Computer Vision and Pattern Recognition*, pp. 324–333, 2019. 3

- 379 Shengyi Gao, Zhe Chen, Guo Chen, Wenhai Wang, and Tong Lu. Avsegformer: Audio-visual  
380 segmentation with transformer. *arXiv preprint arXiv:2307.01146*, 2023. 6, 7
- 381 Mariana-Iuliana Georgescu, Eduardo Fonseca, Radu Tudor Ionescu, Mario Lucic, Cordelia Schmid,  
382 and Anurag Arnab. Audiovisual masked autoencoders. *arXiv preprint arXiv:2212.05922*, 2022.  
383 3
- 384 Yuan Gong, Andrew Rouditchenko, Alexander H Liu, David Harwath, Leonid Karlinsky, Hilde  
385 Kuehne, and James Glass. Contrastive audio-visual masked autoencoder. *arXiv preprint*  
386 *arXiv:2210.07839*, 2022. 3
- 387 Dawei Hao, Yuxin Mao, Bowen He, Xiaodong Han, Yuchao Dai, and Yiran Zhong. Improving  
388 audio-visual segmentation with bidirectional generation. *arXiv preprint arXiv:2308.08288*, 2023.  
389 6
- 390 Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing  
391 Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for  
392 large-scale audio classification. In *2017 IEEE international conference on acoustics, speech and*  
393 *signal processing (icassp)*, pp. 131–135. IEEE, 2017. 17
- 394 Di Hu, Feiping Nie, and Xuelong Li. Deep multimodal clustering for unsupervised audiovisual  
395 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-*  
396 *tion*, pp. 9248–9257, 2019. 3
- 397 Li Hu, Peng Zhang, Bang Zhang, Pan Pan, Yinghui Xu, and Rong Jin. Learning position and tar-  
398 get consistency for memory-based video object segmentation. *arXiv preprint arXiv:2104.04329*,  
399 2021. 15
- 400 Ping Hu, Fabian Caba, Oliver Wang, Zhe Lin, Stan Sclaroff, and Federico Perazzi. Temporally  
401 distributed networks for fast video semantic segmentation. In *Proceedings of the IEEE/CVF*  
402 *Conference on Computer Vision and Pattern Recognition*, pp. 8818–8827, 2020. 15
- 403 Suyog Dutt Jain, Bo Xiong, and Kristen Grauman. Fusionseg: Learning to combine motion and ap-  
404 pearance for fully automatic segmentation of generic objects in videos. In *2017 IEEE conference*  
405 *on computer vision and pattern recognition (CVPR)*, pp. 2117–2126. IEEE, 2017. 15
- 406 Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models  
407 from self-supervised synchronization. *Advances in Neural Information Processing Systems*, 31,  
408 2018. 3
- 409 Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics*  
410 *quarterly*, 2(1-2):83–97, 1955. 6
- 411 Xiaoya Li, Xiaofei Sun, Yuxian Meng, Junjun Liang, Fei Wu, and Jiwei Li. Dice loss for data-  
412 imbalanced nlp tasks. *arXiv preprint arXiv:1911.02855*, 2019. 6
- 413 Yule Li, Jianping Shi, and Dahua Lin. Low-latency video semantic segmentation. In *Proceedings of*  
414 *the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5997–6005, 2018. 15
- 415 Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie.  
416 Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on com-*  
417 *puter vision and pattern recognition*, pp. 2117–2125, 2017a. 5
- 418 Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense  
419 object detection. In *Proceedings of the IEEE international conference on computer vision*, pp.  
420 2980–2988, 2017b. 6
- 421 Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, and Ming-Hsuan Yang. Exploring cross-  
422 video and cross-modality signals for weakly-supervised audio-visual video parsing. *Advances in*  
423 *Neural Information Processing Systems*, 34:11449–11461, 2021. 3
- 424 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
425 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint*  
426 *arXiv:2103.14030*, 2021. 6

- 427 Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin trans-  
428 former. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
429 pp. 3202–3211, 2022. 6
- 430 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
431 *arXiv:1711.05101*, 2017. 7
- 432 Sabarinath Mahadevan, Ali Athar, Aljoša Ošep, Sebastian Hennen, Laura Leal-Taixé, and Bastian  
433 Leibe. Making a case for 3d convolutions for object segmentation in videos. *arXiv preprint*  
434 *arXiv:2008.11516*, 2020. 6
- 435 Yuxin Mao, Jing Zhang, Zhexiong Wan, Yuchao Dai, Aixuan Li, Yunqiu Lv, Xinyu Tian, Deng-Ping  
436 Fan, and Nick Barnes. Transformer transforms salient object detection and camouflaged object  
437 detection. *arXiv preprint arXiv:2104.10127*, 2021. 6
- 438 Shentong Mo and Pedro Morgado. A closer look at weakly-supervised audio-visual source localiza-  
439 tion. *arXiv preprint arXiv:2209.09634*, 2022a. 3
- 440 Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In *Computer Vision–*  
441 *ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings,*  
442 *Part XXXVII*, pp. 218–234. Springer, 2022b. 3
- 443 Shentong Mo and Yapeng Tian. Multi-modal grouping network for weakly-supervised audio-visual  
444 video parsing. In *Advances in Neural Information Processing Systems*, 2022. 3
- 445 Pedro Morgado, Nuno Nvasconcelos, Timothy Langlois, and Oliver Wang. Self-supervised genera-  
446 tion of spatial audio for 360 video. *Advances in neural information processing systems*, 31, 2018.  
447 3
- 448 Pedro Morgado, Yi Li, and Nuno Nvasconcelos. Learning representations from audio-visual spatial  
449 alignment. *Advances in Neural Information Processing Systems*, 33:4733–4744, 2020. 3
- 450 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
451 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-  
452 performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019. 7
- 453 Matthieu Paul, Christoph Mayer, Luc Van Gool, and Radu Timofte. Efficient video semantic seg-  
454 mentation with labels propagation and refinement. In *Proceedings of the IEEE/CVF Winter Con-*  
455 *ference on Applications of Computer Vision*, pp. 2873–2882, 2020. 15
- 456 Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and  
457 Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint*  
458 *arXiv:1704.00675*, 2017. 7
- 459 Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources  
460 localization from coarse to fine. In *Computer Vision–ECCV 2020: 16th European Conference,*  
461 *Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pp. 292–308. Springer, 2020. 3, 6
- 462 Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese.  
463 Generalized intersection over union: A metric and a loss for bounding box regression. In *Pro-*  
464 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 658–666,  
465 2019. 6
- 466 Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to local-  
467 ize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision*  
468 *and Pattern Recognition*, pp. 4358–4366, 2018. 3
- 469 Seonguk Seo, Joon-Young Lee, and Bohyung Han. Urvos: Unified referring video object segmen-  
470 tation network with a large-scale benchmark. In *European Conference on Computer Vision*, pp.  
471 208–223. Springer, 2020. 15
- 472 Hongje Seong, Junhyuk Hyun, and Euntai Kim. Kernelized memory network for video object  
473 segmentation. In *European Conference on Computer Vision*, pp. 629–645. Springer, 2020. 15

- 474 Hongje Seong, Seoung Wug Oh, Joon-Young Lee, Seongwon Lee, Suhyeon Lee, and Euntai Kim.  
475 Hierarchical memory matching network for video object segmentation, 2021. [15](#)
- 476 Guolei Sun, Yun Liu, Henghui Ding, Thomas Probst, and Luc Van Gool. Coarse-to-fine feature min-  
477 ing for video semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer*  
478 *Vision and Pattern Recognition*, pp. 3126–3137, 2022. [15](#)
- 479 Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-supervised  
480 audio-visual video parsing. In *Computer Vision–ECCV 2020: 16th European Conference, Glas-*  
481 *gow, UK, August 23–28, 2020, Proceedings, Part III 16*, pp. 436–454. Springer, 2020. [3](#)
- 482 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*  
483 *learning research*, 9(11), 2008. [9](#)
- 484 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, un-  
485 definedukasz Kaiser, and Illia Polosukhin. Attention is all you need. NIPS’17, pp. 6000–6010,  
486 Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964. [17](#)
- 487 Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo,  
488 and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational*  
489 *Visual Media*, 8(3):415–424, 2022. [6](#)
- 490 Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring  
491 video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
492 *Pattern Recognition*, pp. 4974–4984, 2022. [15](#)
- 493 Yu Wu and Yi Yang. Exploring heterogeneous clues for weakly-supervised audio-visual video pars-  
494 ing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
495 pp. 1326–1335, 2021. [3](#)
- 496 Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos K Katsaggelos. Efficient  
497 video object segmentation via network modulation. In *Proceedings of the IEEE Conference on*  
498 *Computer Vision and Pattern Recognition*, pp. 6499–6507, 2018. [15](#)
- 499 Zongxin Yang, Yunchao Wei, and Yi Yang. Associating objects with transformers for video object  
500 segmentation. *Advances in Neural Information Processing Systems*, 34, 2021. [6](#), [7](#), [15](#)
- 501 Jing Zhang, Jianwen Xie, Nick Barnes, and Ping Li. Learning generative vision transformer with  
502 energy-based latent space for saliency prediction. *Advances in Neural Information Processing*  
503 *Systems*, 34:15448–15463, 2021. [6](#)
- 504 Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio  
505 Torralba. The sound of pixels. In *Proceedings of the European conference on computer vision*  
506 *(ECCV)*, pp. 570–586, 2018. [3](#)
- 507 Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In *Pro-*  
508 *ceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1735–1744, 2019.  
509 [3](#)
- 510 Jinxing Zhou, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birchfield, Dan Guo,  
511 Lingpeng Kong, Meng Wang, and Yiran Zhong. Audio-visual segmentation. In *European Con-*  
512 *ference on Computer Vision*, 2022. [1](#), [2](#), [3](#), [7](#), [8](#)
- 513 Jinxing Zhou, Xuyang Shen, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birch-  
514 field, Dan Guo, Lingpeng Kong, Meng Wang, et al. Audio-visual segmentation with semantics.  
515 *arXiv preprint arXiv:2301.13190*, 2023. [1](#), [2](#), [3](#), [6](#), [7](#), [16](#), [18](#), [19](#), [20](#), [21](#)
- 516 Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:  
517 Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020.  
518 [17](#)
- 519 Jiafan Zhuang, Zilei Wang, and Yuan Gao. Semi-supervised video semantic segmentation with inter-  
520 frame feature reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
521 *and Pattern Recognition*, pp. 3263–3271, 2022. [15](#)



## 522 A MORE EXPERIMENTS

frame number	Object-M	Sementic
	$\mathcal{J}\&\mathcal{F}$	mIoU
3	60.9	45.4
5	61.6	46.6
7	-	46.6

Table 6: Ablation on the input frame number.

523 **Frame number.** We ablate the influence of input frame number during training. As shown in  
 524 Table 6, we notice a frame number of five achieves the best performance. For the AVS-Object  
 525 dataset, since the maximum clip length is five, we do not experiment with larger frame number.  
 526 Please note that the frame number is only fixed during training and the model can accept arbitrary  
 527 frame numbers during inference.

layer number	Object-M	Sementic
	$\mathcal{J}\&\mathcal{F}$	mIoU
1	61.3	45.6
3	61.6	46.6
5	61.0	46.0

Table 7: Ablation on transformer decoder layer number.

528 **Transformer decoder layer number.** We conduct an ablation study on transformer decoder layer  
 529 numbers in semantic decoders. As shown in Table 7, a transformer decoder layer of 3 achieves the  
 530 best performance. We notice that even a single-layer transformer decoder for semantic decomposi-  
 531 tion can lead to a good performance.

frame resolution	Sementic
	mIoU
224×	46.6
640×	49.2

Table 8: Ablation on input frame resolution.

532 **Input Resolution.** The default setting of AVSBench is  $224 \times 224$  (following the sound source local-  
 533 ization convention) for both AVS-Object and AVS-Semantic datasets. While AVS-Semantic actually  
 534 provides high-resolution (720p) frames. We conduct experiments to ablate the input resolution to  
 535 facilitate future comparison. Following the semantic segmentation convention, we scale the input  
 536 frames to the longest side 224 or 640. The results are illustrated in Table 8. We only conduct ablation  
 537 on AVS-Semantic since the resolution of AVS-Object is low-resolution ( $224 \times 224$ ). The results are  
 538 reported with the ResNet-50 backbone.

539 **Per-class IOU analysis.** As is shown in Fig. 8, we show the per-class iou score on the AVS-Semantic  
 540 dataset. Our model demonstrates strong audio-guided segmentation capabilities for common head  
 541 classes such as 'background', 'train', 'airplane', 'hair-dryer' and 'clock'. These classes are accu-  
 542 rately segmented with a high level of precision and reliability. The model effectively distinguishes  
 543 the 'background' class, providing a solid foundation for identifying and isolating foreground objects.  
 544 It accurately segments transportation-related classes like 'train', 'airplane', and 'bus' capturing their  
 545 intricate details and boundaries. Similarly, it excels in segmenting objects such as 'hair-dryer',  
 546 'clock' and 'tabla,' effectively separating them from the background. Even for more complex and  
 547 nuanced classes like 'wolf,' our model demonstrates commendable segmentation performance, ac-  
 548 curately delineating the contours and shape of the subject. Overall, our model showcases its ability  
 549 to segment these common head classes with high accuracy and proficiency, making it a reliable  
 550 choice for various segmentation tasks.

551 However, the scarcity of data samples for tail classes like 'utv', 'parrot', 'missile-rocket', 'harmon-  
 552 ica', 'clipper', 'boy' and 'ax' in the presence of a long tail distribution can significantly impact the

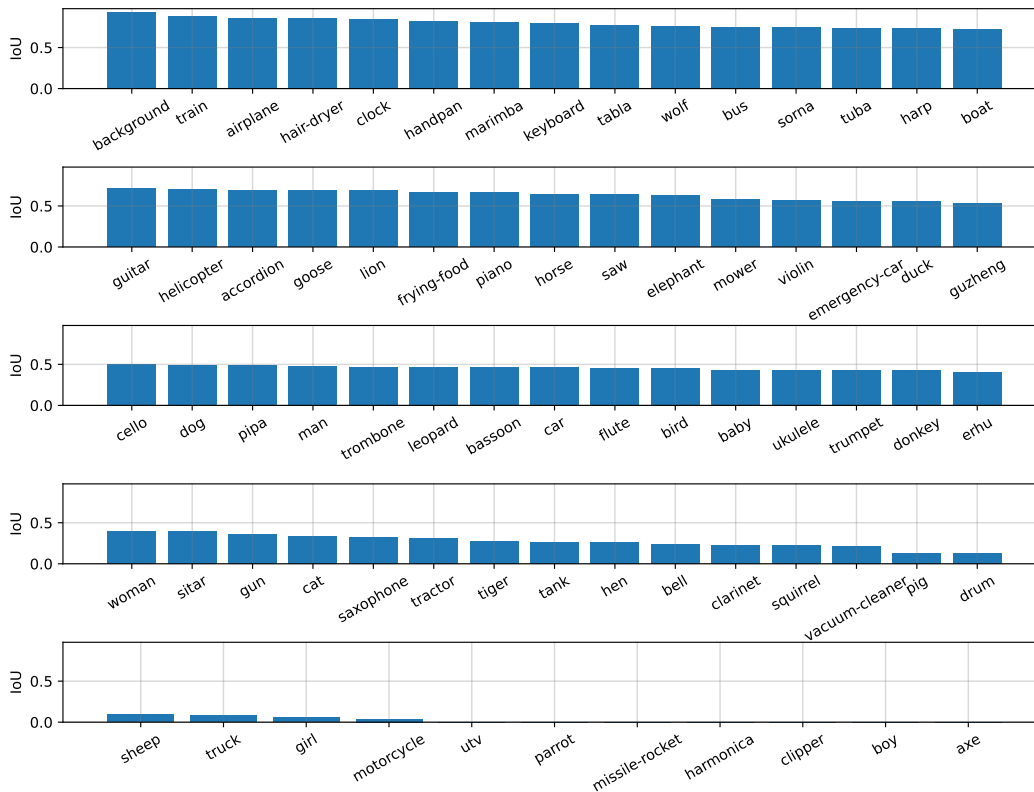


Figure 8: **Per-class IOU Analysis.** Our model demonstrates strong audio-guided segmentation capabilities for common head classes, accurately capturing 'background', 'train', 'airplane', 'hair-dryer', and 'clock' with high precision. However, the limited data samples for tail classes like 'utv', 'parrot', 'missile-rocket', 'harmonica', 'clipper', 'boy', and 'ax' due to a long tail distribution adversely affect the model's segmentation performance, hindering accurate identification and delineation of these classes.

553 performance of our model, specifically in the task of segmentation. With limited examples to learn  
 554 from, the model finds it challenging to capture the intricate patterns and unique characteristics asso-  
 555 ciated with these classes. Consequently, the accuracy and reliability of segmentation results for the  
 556 tail classes may be compromised, leading to suboptimal performance in accurately identifying and  
 557 delineating these objects or entities of interest.

## 558 B MORE RELATED WORKS

559 Audiovisual segmentation also closely relates to video object segmentation (VOS) Yang et al.  
 560 (2018); Jain et al. (2017); Cheng et al. (2021b); Seong et al. (2020); Hu et al. (2021); Cheng et al.  
 561 (2021c); Seong et al. (2021); Yang et al. (2021) and video semantic segmentation (VSS) Li et al.  
 562 (2018); Sun et al. (2022); Zhuang et al. (2022); Hu et al. (2020); Paul et al. (2020). AVS requires  
 563 understanding the visual contents and then corresponding them with the audio semantics to segment  
 564 objects. Specifically, the most closely related task in the video segmentation domain is the referring  
 565 video object segmentation (RVOS) Botach et al. (2022); Wu et al. (2022); Seo et al. (2020) which  
 566 aims to segment objects in the visual frames given a linguistic expression. For each expression,  
 567 RVOS only refers to one object while the AVS task permits audio query to refer to multiple objects  
 568 which makes AVS task more challenging.

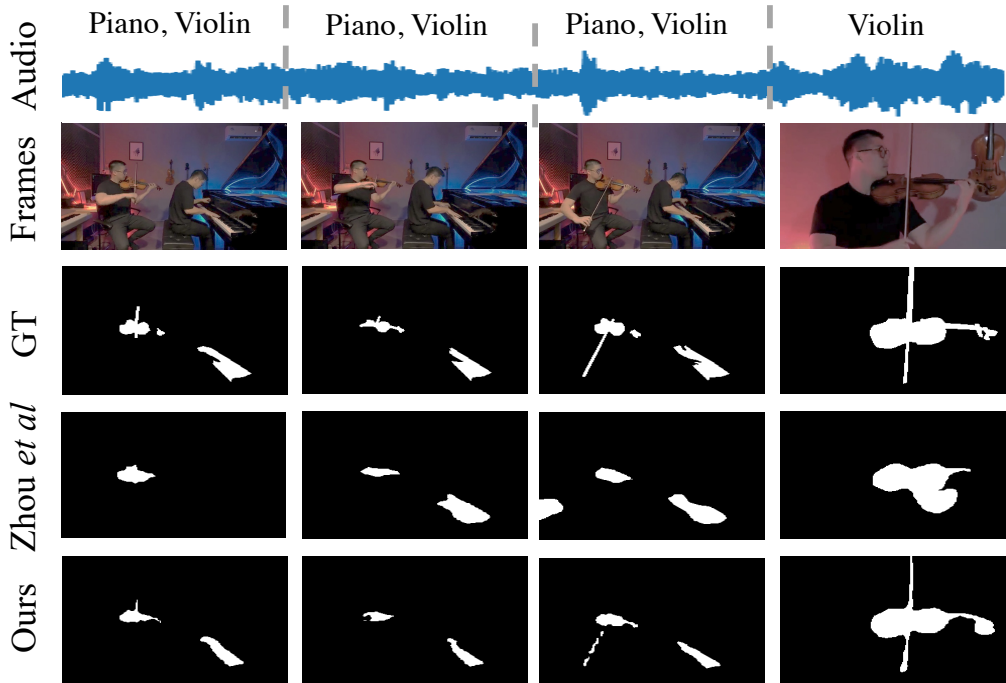


Figure 9: Qualitative comparison to Zhou et al. [Zhou et al. \(2023\)](#) on AVS-Object. Our method outperforms Zhou et al.’s approach by consistently and accurately segmenting the correct objects throughout the entire video clip, showcasing superior performance and better mask quality. These results emphasize the effectiveness and robustness of our approach in achieving accurate object segmentation in audio-visual scenes.

569 C MORE VISUALIZATION & VIDEO DEMO

570 **More qualitative results on AVS-Object.** In our study, we provide visualizations of the qualitative  
 571 results on AVS-Object, as shown in Fig. 9. We compare our method with the approach proposed  
 572 by Zhou et al. [Zhou et al. \(2023\)](#) and observe a notable difference in performance. Specifically, in  
 573 the third frame of the video clip, the method proposed by Zhou et al. suffers from the false-positive  
 574 problem, incorrectly segmenting objects. In contrast, our method consistently and accurately seg-  
 575 ments the correct objects throughout the entire video clip, demonstrating superior performance. Ad-  
 576 ditionally, our method showcases better mask quality, with more precise and detailed segmentation  
 577 boundaries. These results highlight the effectiveness and robustness of our approach in achieving  
 578 accurate object segmentation in audio-visual scenes.

579 **More qualitative results on AVS-Semantics.** As is shown in Fig. 10, Fig. 11, Fig. 12 and Fig. 13,  
 580 our model exhibits exceptional proficiency in accurately segmenting both multiple and tiny sound-  
 581 ing objects, showcasing its versatility and robustness in audio-guided segmentation tasks. Through  
 582 the implementation of a decomposed and discretized audio representation, our model effectively  
 583 captures the distinct acoustic characteristics of various objects, enabling precise delineation of mul-  
 584 tiple simultaneous sound sources. Furthermore, the model demonstrates remarkable capability in  
 585 capturing the intricate details and nuances of tiny sounding objects, ensuring accurate segmentation  
 586 outcomes even for the smallest entities.

587 **Video demo (with audio).** We strongly recommend viewing the demo video provided in the sup-  
 588 plementary materials, ensuring that you enable audio playback. Watching the video with audio will  
 589 provide a comprehensive understanding of our audio-visual segmentation application, showcasing  
 590 how our model utilizes a decomposed and discretized representation to achieve precise audio-visual  
 591 segmentation results.

## 592 D MORE IMPLEMENTATION DETAILS

593 We set the  $\lambda_{cls} = 2$ ,  $\lambda_{L1} = 5$ ,  $\lambda_{giou} = 2$ ,  $\lambda_{dice} = 2$ ,  $\lambda_{focal} = 5$ ,  $\lambda_{com} = 0.5$  and  $\lambda_{quant} = 1$  during  
 594 all training process. A mask confidence threshold of 0.5 and a class confidence threshold of 0.1 is  
 595 leveraged to filter out low-confident predictions.  $C_v = C_e = C_q = 256$  is utilized. The positional  
 596 embedding added in the transformers is the standard triangle positional embedding used in [Vaswani](#)  
 597 [et al. \(2017\)](#). We set the layer number to three for all the transformers decoders (including local  
 598 ASD, global ASD and TrD<sub>segm</sub> in mask decoder).

### 599 D.1 ENCODERS

600 **Visual encoder.** The visual encoder consists of a visual backbone and a deformable transformer  
 601 encoder [Zhu et al. \(2020\)](#). We extract frame-level visual features from each frame  $I_t$  with a shared  
 602 backbone. The  $T$  extracted features are then fed into the deformable transformer encoder to further  
 603 conduct temporal aggregation. Let us denote the extracted visual features as  $F_v = \{f_t\}_{t=1}^T$ , where  
 604  $f_t \in \mathbb{R}^{C_v \times H \times W}$ , and  $C_v, H, W$  denote the channel, height, width of the feature.

605 **Acoustic encoder.** We use VGGish [Hershey et al. \(2017\)](#) to extract audio features. Let the extracted  
 606 audio feature be  $F_a \in \mathbb{R}^{C_a \times L_a}$  where  $C_a$  is the dimension of acoustic feature space, and  $L_a$  is the  
 607 audio clip length. Note that audio and video frames are already synchronized, thus the length of the  
 608 audio clip is the same as the length of the video clip.

## 609 E MORE DETAILS ABOUT INFERENCE

610 To tackle scenarios where queried content keeps changing, we perform per-frame inference. For  
 611 each time  $t$ , we assign a class to the pixel at  $[h, w]$  by

$$\arg \max_{C \in \{1, \dots, K\}} \sum_{i=1}^N P_{i,t}[C] M_{i,t}[h, w], \quad (11)$$

612 where  $P_{i,t}[C]$  is the probability of class  $C$ . Note that  $\arg \max$  does not include the “empty” category  
 613 ( $\emptyset$ ) as AVS requires each output pixel to belong to one semantic category.

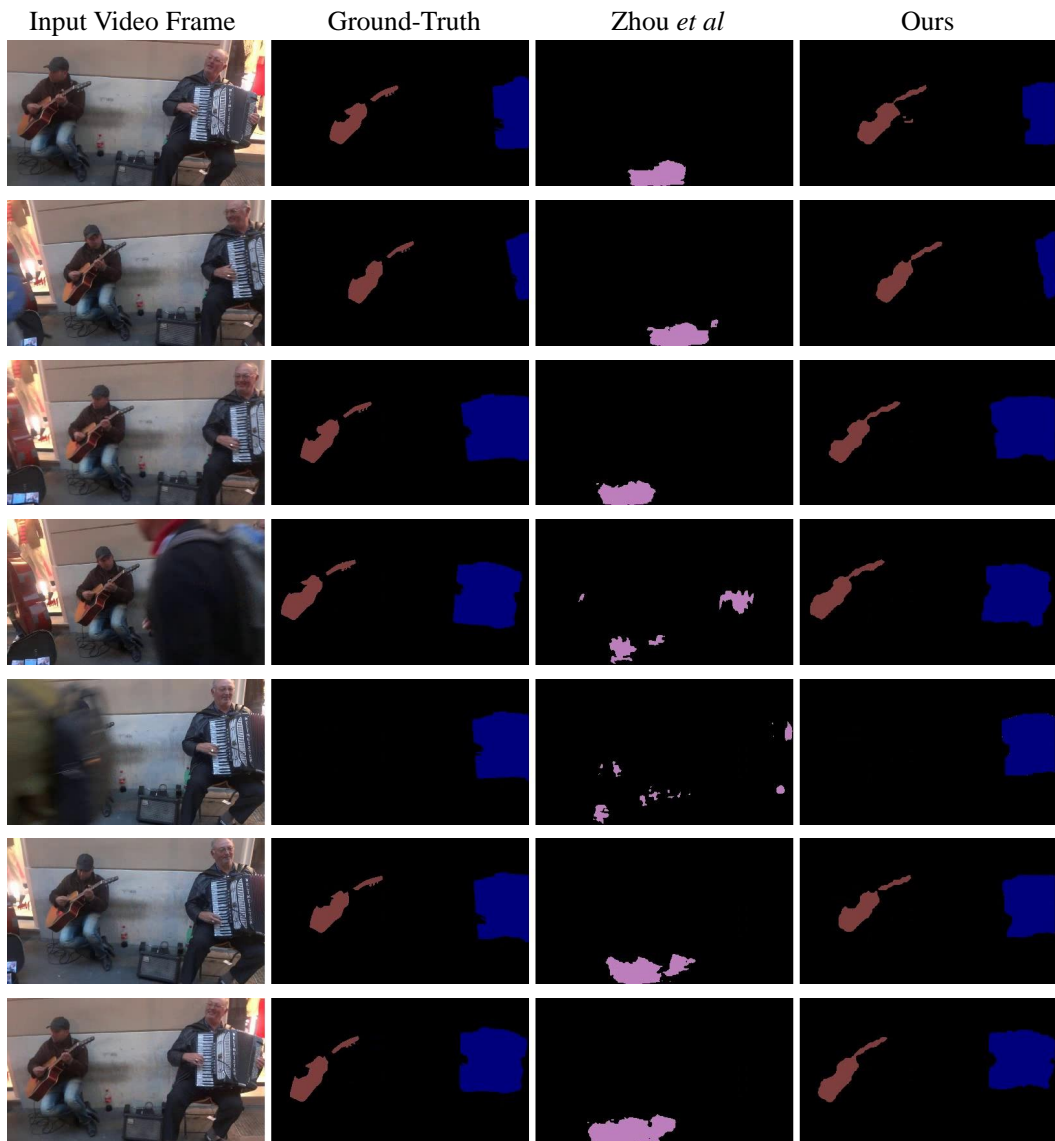


Figure 10: Qualitative comparison to Zhou et al. [Zhou et al. \(2023\)](#) on AVS-Semantic. Each color represents a semantic category. Our model excels in accurately segmenting **multiple sounding objects**, showcasing its proficiency in audio-guided segmentation. This success can be attributed to the effective utilization of a decomposed and discretized audio representation, which enables the model to capture and analyze the distinct acoustic features of each object, resulting in precise segmentation outcomes.



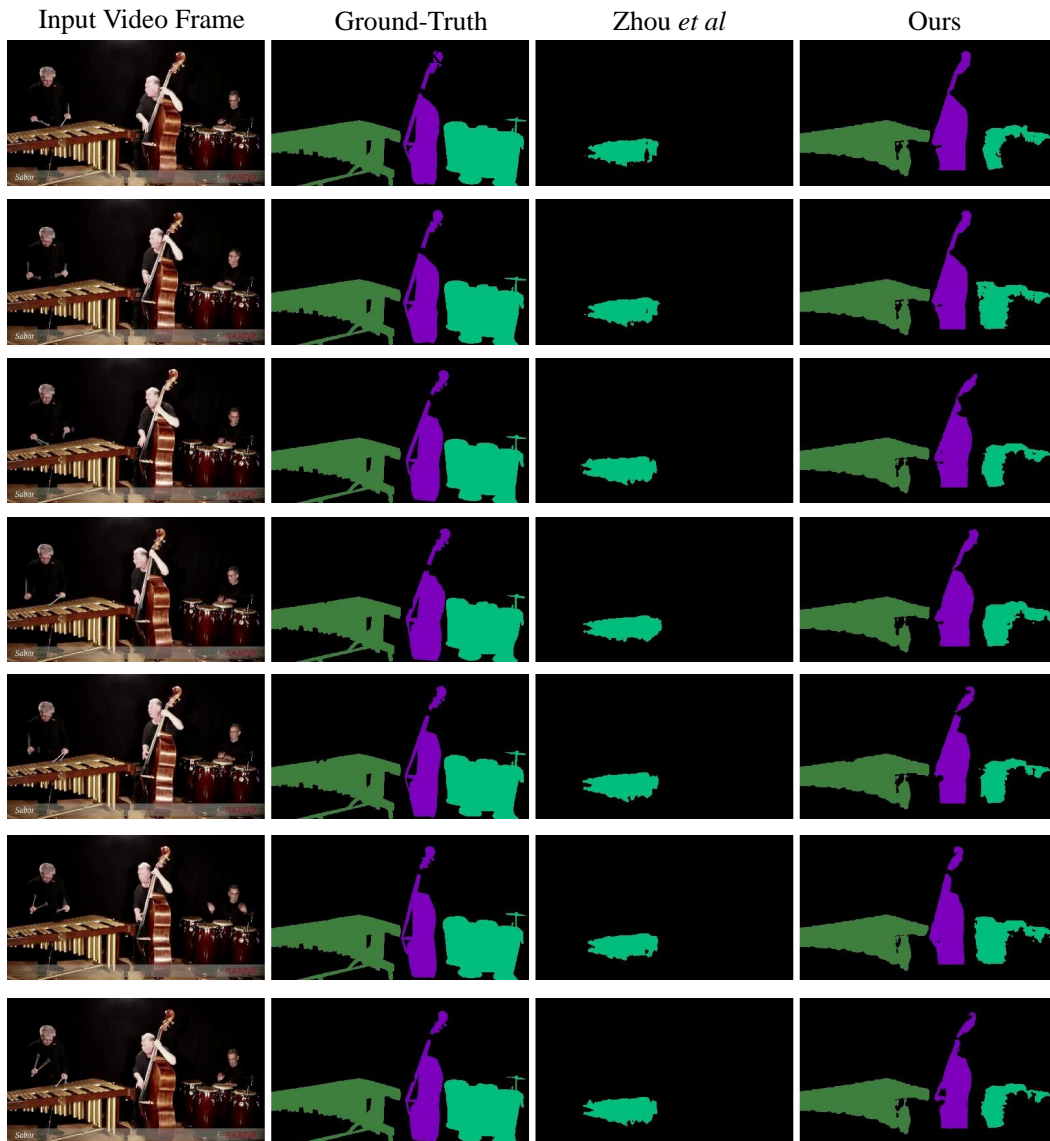


Figure 11: Qualitative comparison to Zhou et al. [Zhou et al. \(2023\)](#) on AVS-Semantic. Each color represents a semantic category. Our model excels in accurately segmenting **multiple sounding objects**, showcasing its proficiency in audio-guided segmentation. This success can be attributed to the effective utilization of a decomposed and discretized audio representation, which enables the model to capture and analyze the distinct acoustic features of each object, resulting in precise segmentation outcomes.



Figure 12: Qualitative comparison to Zhou et al. [Zhou et al. \(2023\)](#) on AVS-Semantic. Each color represents a semantic category. Our model excels in accurately segmenting **multiple sounding objects**, showcasing its proficiency in audio-guided segmentation. This success can be attributed to the effective utilization of a decomposed and discretized audio representation, which enables the model to capture and analyze the distinct acoustic features of each object, resulting in precise segmentation outcomes.

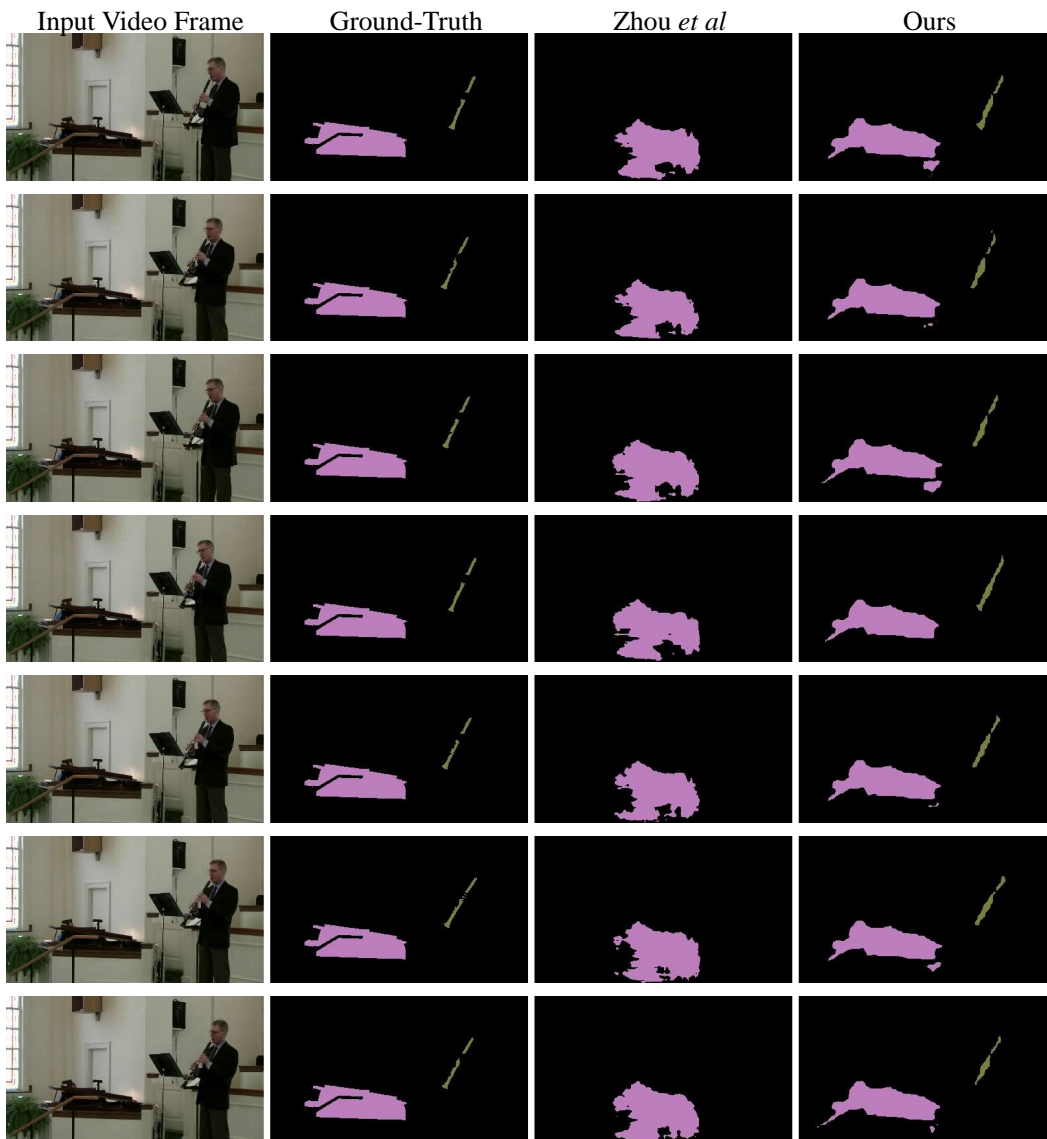


Figure 13: Qualitative comparison to Zhou et al. [Zhou et al. \(2023\)](#) on AVS-Semantic. Each color represents a semantic category. Our model demonstrates remarkable proficiency in accurately segmenting **tiny sounding objects**, owing to the implementation of a decomposed and discretized audio representation. By leveraging this technique, our model effectively captures the intricate acoustic details and nuances of these small-sized objects, resulting in precise and reliable segmentation outcomes.