



\*Jinxing Zhou<sup>1,2</sup>, \*Jianyuan Wang<sup>3</sup>, Jiayi Zhang<sup>2,4</sup>, Weixuan Sun<sup>2,3</sup>, Jing Zhang<sup>3</sup>, Stan Birchfield<sup>5</sup>, Dan Guo<sup>1</sup>, Lingpeng Kong<sup>6,7</sup>, Meng Wang<sup>1</sup>, Yiran Zhong<sup>2,7</sup>

<sup>1</sup>Hefei University of Technology, <sup>2</sup>SenseTime Research, <sup>3</sup>Australian National University, <sup>4</sup>Beihang University, <sup>5</sup>NVIDIA, <sup>6</sup>The University of Hong Kong, <sup>7</sup>Shanghai Artificial Intelligence Laboratory



## Introduction

### Motivation

- Recent research on audio-visual learning mainly focuses on coarse-grained scene understanding, such as the audio-visual event classification and the patch-level sound source localization (SSL).
- We propose the audio-visual segmentation (AVS) task, which aims to explore the pixel-level audio-visual correspondence. The goal of AVS is to output a pixel-level segmentation map of the object(s) that produce sound at the time of the image frame (Figure 1).

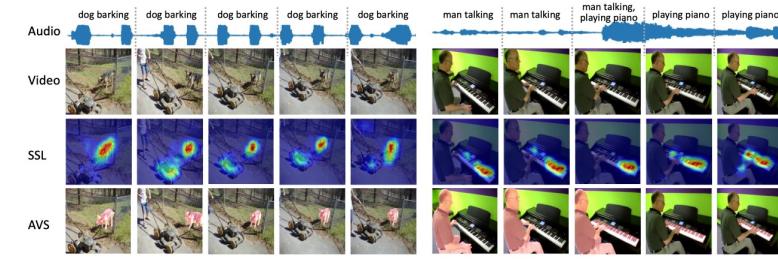


Figure 1. Comparison of the proposed AVS task with the SSL task.

### AVSBench Dataset

To facilitate the research of AVS, we collect the AVSBench, a new dataset providing pixel-level annotations.

- Data Statistics.** AVSBench contains 5,356 videos in total covering 23 types of common audio-visual scenes in real life. AVSBench can be divided into two subsets, i.e., the Single-source and the Multi-sources subset, according to the number of sound sources (Table 1).

subset	classes	videos	train/valid/test	annotated frames
Single-source	23	4,932	3,452*/740/740	10,852
Multi-sources	23	424	296/64/64	2,120

Table 1. AVSBench statistics.

- Annotation.** For the Single-source subset, only the first sampled frame is annotated for the training process. As for the Multi-sources subset, all of the five sampled frames are annotated (Figure 2). This induces the two settings of the AVS task, i.e., the semi-supervised Single Sound Source Segmentation (S4) and the fully-supervised Multiple Sound Source Segmentation (MS3).



Figure 2. AVSBench samples.

## Audio-Visual Segmentation

### Baseline Method

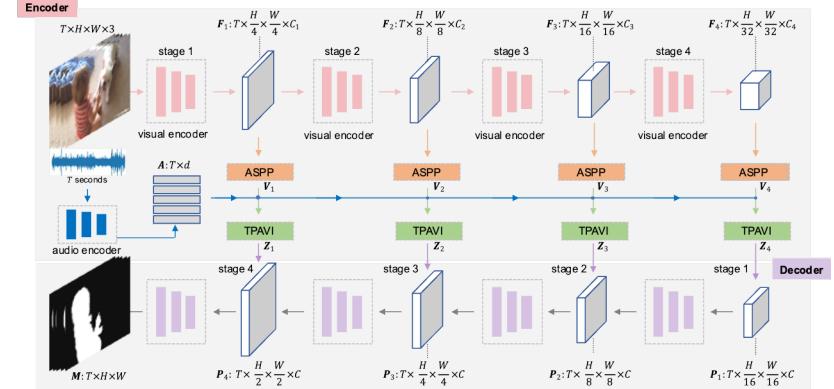


Figure 3. Overview of the Baseline Framework.

### Baseline Encoder-Decoder Framework.

- The *encoder* takes the video frames and the entire audio clip as input, and outputs the visual and audio features, respectively denoted as  $\mathbf{F}_i$  and  $\mathbf{A}$ . The visual feature map  $\mathbf{F}_i$  at each stage is further sent to the ASPP [7] module and then our TPAVI module (Figure 4) which focuses on the temporal pixel-wise audio-visual interaction.
- The *decoder* progressively enlarges the fused feature maps by four stages and finally generates the output mask  $\mathbf{M}$  for sounding objects.

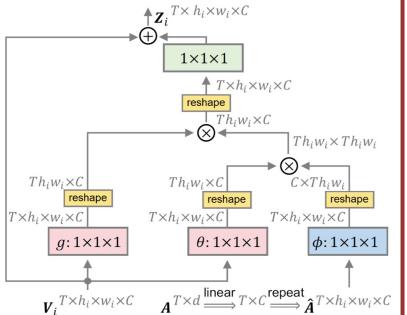


Figure 4. Illustration of the TPAVI module.

**Objective Function.** Given the prediction  $\mathbf{M}$  and the pixel-wise label  $\mathbf{Y}$ , we adopt the binary cross entropy (BCE) loss as the main supervision functions.

Besides, we use an additional regularization term  $\mathcal{L}_{AVM}$  to force the audio-visual mapping (Eq. 2).  $\mathcal{L}_{AVM}$  computes the Kullback-Leibler (KL) divergence to ensure the masked visual features have similar distributions with the corresponding audio features.

$$\mathcal{L} = \text{BCE}(\mathbf{M}, \mathbf{Y}) + \lambda \mathcal{L}_{AVM}(\mathbf{M}, \mathbf{Z}, \mathbf{A}), \quad (1)$$

$$\mathcal{L}_{AVM} = \sum_{i=1}^n (\text{KL}(\text{avg}(\mathbf{M}_i \odot \mathbf{Z}_i), \mathbf{A}_i)), \quad (2)$$

### Experiments

Comparison with methods from related tasks.

Metric	Setting	SSL		VOS		SOD		AVS (ours)	
		LVS [5]	MSSL [30]	3DC [27]	SST [10]	iGAN [28]	LGVT [19]	ResNet50	PVT-v2
$\mathcal{M}_{\mathcal{I}}$	S4	.379	.449	.571	.663	.616	.749	.728	<b>.787</b>
	MS3	.295	.261	.369	.426	.429	.407	.479	<b>.540</b>
$\mathcal{M}_{\mathcal{F}}$	S4	.510	.663	.759	.801	.778	.873	.848	<b>.879</b>
	MS3	.330	.363	.503	.572	.544	.593	.578	<b>.645</b>

Impact of audio signal and the proposed TPAVI module.

AVS method	S4		MS3	
	ResNet50	PVT-v2	ResNet50	PVT-v2
without TPAVI	.701	.778	.436	.482
with A $\oplus$ V	.705	.777	.457	.516
with TPAVI	<b>.728</b>	<b>.787</b>	<b>.466</b>	<b>.531</b>

Effectiveness of the proposed  $\mathcal{L}_{AVM}$ .

Objective function	MS3 (mIoU)		MS3 (F-score)	
	ResNet50	PVT-v2	ResNet50	PVT-v2
$\mathcal{L}_{BCE}$	46.64	53.06	.558	.626
$\mathcal{L}_{BCE} + \mathcal{L}_{AVM-VV}$	46.71	53.77	.577	.644
$\mathcal{L}_{BCE} + \mathcal{L}_{AVM-AV}$	<b>47.88</b>	<b>54.00</b>	<b>.578</b>	<b>.645</b>

Qualitative results under the semi-supervised S4 setting.



Qualitative results under the fully-supervised MS3 setting.

