

Audio-Visual Segmentation by Exploring Cross-Modal Mutual Semantics

Chen Liu^{1,2,5}, Peike Li³, Xingqun Qi^{1,4}, Hu Zhang², Lincheng Li⁶, Dadong Wang⁵, Xin Yu²

¹University of Technology Sydney ²The University of Queensland ³Futureverse

⁴The Hong Kong University of Science and Technology ⁵CSIRO DATA61 ⁶Netease Fuxi AI Lab

{yenanliu36,xingqunqi}@gmail.com,{xin.yu,hu.zhang}@uq.edu.au

peike.li@yahoo.com,lilincheng@corp.netease.com,Dadong.Wang@data61.csiro.au

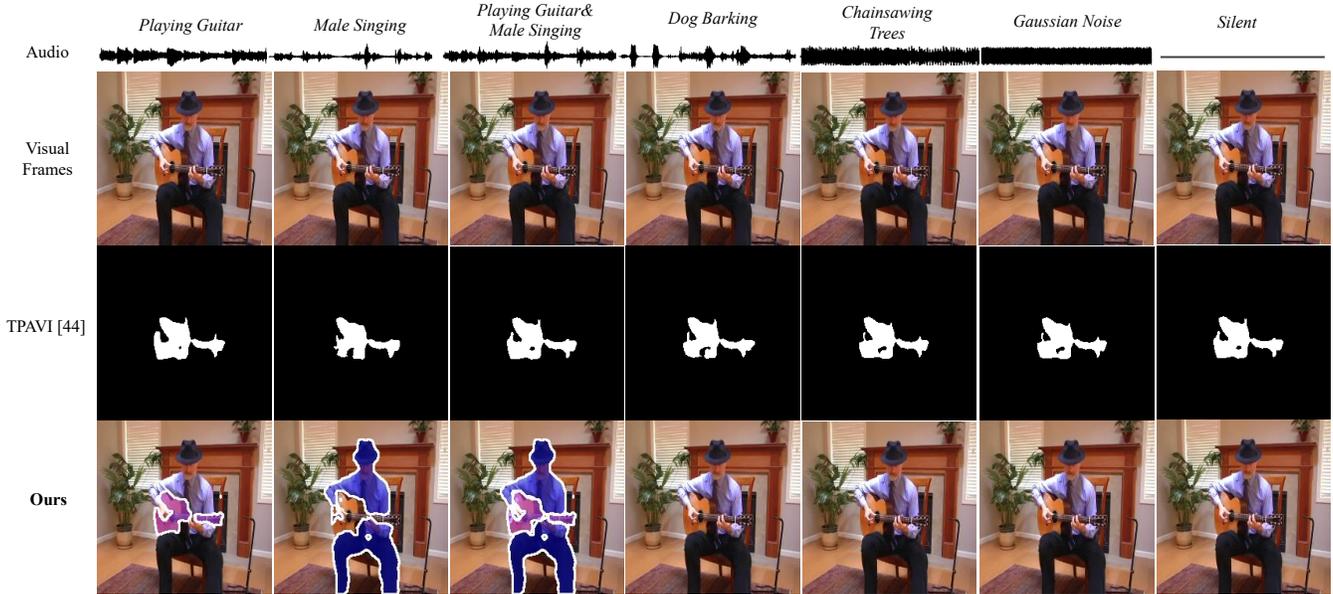


Figure 1: Visual results of our method and TPAVI [46]. Our method accurately segments sounding objects according to different audio signals. For the audio signal mixed by "Male Singing" and "Playing Guitar", our method even achieves instance-level sounding object segmentation. Here, we employ different colored masks to represent segmentation masks of different instances. Besides, our method is not misled when the guitar signal is muted or replaced by "Dog Barking", "Chainsawing Trees", or "Gaussian Noise". In contrast, the audio semantic information fails to modulate segmentation results in [46]. No matter what the audio signal is, TPAVI always segments the most prominent guitar.

ABSTRACT

The audio-visual segmentation (AVS) task aims to segment sounding objects from a given video. Existing works mainly focus on fusing audio and visual features of a given video to achieve sounding object masks. However, we observed that prior arts are prone to segment a certain salient object in a video regardless of the audio information. This is because sounding objects are often the most salient ones in the AVS dataset. Thus, current AVS methods might fail to localize genuine sounding objects due to the dataset bias. In this work, we present an audio-visual instance-aware segmentation approach to overcome the dataset bias. In a nutshell, our method first localizes potential sounding objects in a video by an object segmentation network, and then associates the sounding object candidates with the given audio. We notice that an object could be a sounding object in one video but a silent one in another video. This would bring ambiguity in training our object segmentation network

as only sounding objects have corresponding segmentation masks. We thus propose a silent object-aware segmentation objective to alleviate the ambiguity. Moreover, since the category information of audio is unknown, especially for multiple sounding sources, we propose to explore the audio-visual semantic correlation and then associate audio with potential objects. Specifically, we attend predicted audio category scores to potential instance masks and these scores will highlight corresponding sounding instances while suppressing inaudible ones. When we enforce the attended instance masks to resemble the ground-truth mask, we are able to establish audio-visual semantics correlation. Experimental results on the AVS benchmarks demonstrate that our method can effectively segment sounding objects without being biased to salient objects and also achieves state-of-the-art performance in both the single-source and multi-source scenarios.

KEYWORDS

Audio-visual segmentation, sound localization, semantic-aware sounding objects localization

1 INTRODUCTION

Sounding source localization has been applied in various multimedia applications, such as robotic navigation [40], security monitoring [41], wildlife conservation [30], and industrial maintenance [13]. The Audio-Visual Segmentation (AVS) task [46] further pushes the frontier of sounding source localization and strives to localize the masks of sounding objects in a pixel-wise manner.

Existing audio-visual segmentation methods predominantly first fuse audio-visual features and then predict sounding object masks. However, we found that current methods are prone to segment a certain salient object in a video regardless of audio information. For instance, as illustrated in Figure 1, when we replace the original audio (*i.e.*, playing guitar) with the audio of another object (*i.e.*, male singing), the segmentation result of [46] remains almost the same. Furthermore, when the audio is muted or replaced by “dog barking”, a blank mask is expected but [46] still outputs the mask of the guitar. This implies that audio information might not play its role effectively in existing AVS approaches. As a sounding object is often the most salient object in a video, this bias leads current AVS networks to localize a specific salient object while neglecting other potential audible objects. Therefore, current AVS methods would fail to segment sounding objects in accordance with different audio information.

To tackle the aforementioned issues, we introduce an audio-visual instance-aware segmentation approach. Unlike previous works that focus on fusing audio and visual features and then predicting segmentation masks, our method firstly attempts to localize all the potential sounding objects by learning an instance segmentation neural network. However, we notice that an object could be a sounding object in one video but is an inaudible one in another video but the AVS dataset only provides the segmentation masks of sounding objects. This would lead to ambiguity in learning the segmentation network, and the segmentation network may fail to localize all potential-sounding object masks. To mitigate this issue, we propose a silent object-aware segmentation loss. It not only supervises sounding object segmentation but also can tolerate other object masks which might be potential sounding objects in the dataset. In this way, we can segment all potential sounding objects from a video.

After obtaining sounding object candidates, we seek to associate the audio with genuine sounding instances. Note that the category information of audio is unknown, especially for multiple-sounding sources. Thus, we develop an audio-visual semantic correlation (AVSC) module to mine the audio and visual correspondences. To be specific, our AVSC first estimates category probabilities of audio via some multi-layer perceptron (MLP) layers, and then attends estimated probabilities to their corresponding candidate masks to produce a sounding object probability map. In learning AVSC, we force the sounding object probability map to be similar to the ground-truth mask. If an object instance is a sounding one, its audio category probability will be encouraged to be high. On the other

hand, if an object is inaudible, its estimated audio category probability will be restrained. In this fashion, we effectively associate audio information with sounding objects while significantly alleviating over-fitting to certain salient objects.

Extensive experiments on the widely-used AVS benchmark [46] demonstrate that our method achieves state-of-the-art performance on both single-source and multi-source audio-visual segmentation tasks. More importantly, we also illustrate that our segmentation results can be controlled by various audio signals, demonstrating our method is aware of audio information in segmenting sounding objects. In summary, our contributions are threefold:

- We first present an instance-aware audio-visual segmentation method that allows the segmentation results to be controlled by different audio signals, rather than overfitting to certain salient objects in the dataset.
- We introduce a silent object-aware segmentation objective to alleviate the segmentation supervision provided by the existing AVS dataset. With this objective, we can effectively segment potential sounding instances.
- We design an audio-visual semantic correlation (AVSC) module to establish the association between audio and visual information. Our proposed AVSC enables us to identify the sounding instances from all the candidates without knowing the category information of audio sources.

2 RELATED WORKS

Visual Sound Localization Sounding source localization methods aim to locate the regions in the visual frames related to the corresponding audio signals. When the pixel-level annotations of sounding object locations are not available, previous methods mainly employ coarse heat maps or bounding boxes to localize sounding objects [2, 14, 19, 23–25, 28, 31–34, 36, 42, 45, 47, 48]. Those methods mainly first establish correspondences between audio and visual representations through contrastive learning [2, 23, 24, 32, 36] and then localize objects based on the audio-visual feature similarity. For example, [2] employs contrastive learning with hard negative mining to learn co-occurrence between audio and images, while [32] proposes to leverage hard positive samples to align visual and audio features. The work [28] further develops a multi-instance contrastive learning fashion to match audio to video frames. However, these methods usually localize visible sound sources but struggle to identify negative cases, where sounding objects are not visible. Apart from the difficulty of identifying invisible sounding objects, several methods [14, 19] require knowing the number of sound sources in advance. However, this requirement might be not met, thus limiting their applications in practice.

Since previous methods cannot localize sounding objects accurately, [19, 46] first proposes an audio-visual segmentation dataset with pixel-level audio-visual annotation, allowing researchers to localize sounding objects more precisely. Our work also focuses on predicting sounding object masks while aiming to achieve the ability to recognize invisible sounding sources and fully leverage audio information in segmentation.

Guided Segmentation Networks with Conditional Information. Semantic or instance segmentation tasks play critical roles in

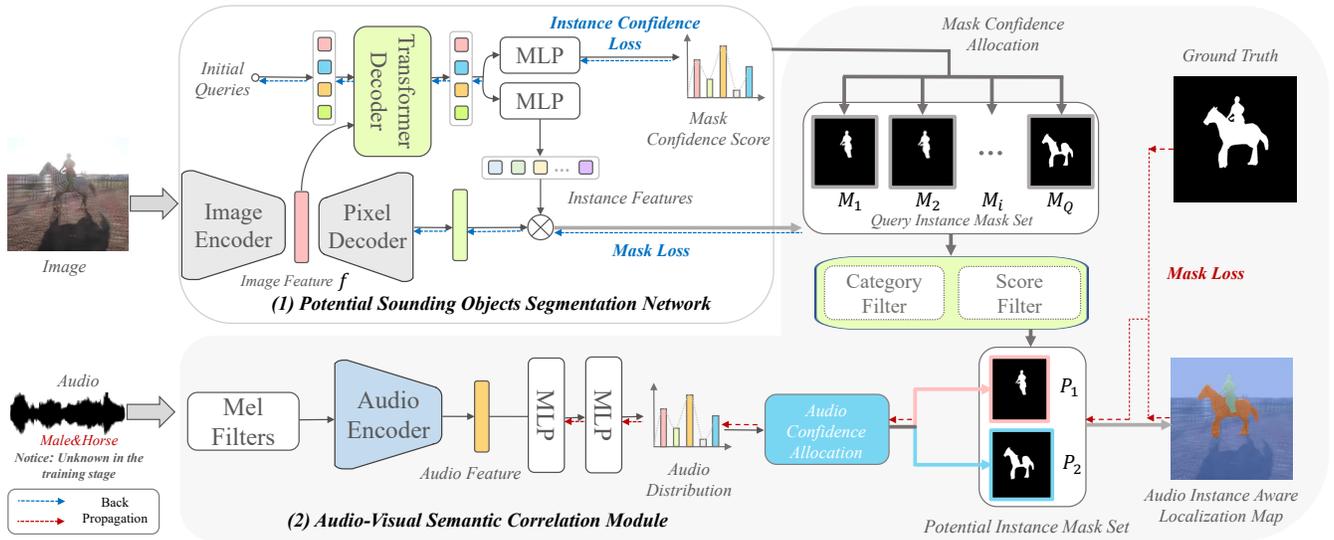


Figure 2: Overview of our audio-visual instance-aware segmentation method. It consists of a potential sounding object segmentation network and an audio-visual semantic correlation module (AVSC). The segmentation network localizes potential-sounding objects. Then AVSC associates the audio with potential sounding instances by attending our estimated audio category labels to instance segmentation masks. During training, sounding objects’ masks will be highlighted while silent ones will be suppressed.

computer vision. They both aim to identify objects and their boundaries within an image [5, 43]. Convolutional neural network (CNN) models have been widely used for these tasks with popular architectures, such as U-Net [29], Mask R-CNN [12], and DeepLab [3, 4]. Recently, transformer-based models have also shown promising results for segmentation tasks. Transformers [38], originally proposed for natural language processing, have shown effectiveness in capturing long-range dependencies and contextual information in images [26, 27]. Some examples of transformer-based segmentation models include DETR [1], SETR [37], and MaskFormer [8]. In addition to these model architectures, researchers have also investigated using different contextual information, as additional conditions, to improve segmentation performance. For instance, text [21] can be used to specify the desired output segmentation masks, depth [39] or surface normal maps [10] can be used as input features, and audio [31, 42] provides additional contextual information to localize specific sounding instances.

Audio-guided segmentation aims to leverage audio signals to guide the segmentation process. Previous works design various mechanisms [18, 34, 35, 46] to combine information from audio and visual information and then perform segmentation from the fused audio-visual features. Although the audio-guided segmentation method has shown promising results [46], it cannot distinguish whether the sound comes from one object or multiple ones. Thus, they cannot further divide a semantic segmentation mask into instance-level masks. Furthermore, due to the bias of the current dataset [46], audio-guided segmentation might favor segmenting salient objects while ignoring audio information. Thus, audio cannot play its role during segmentation. In contrast, the segmentation results of our proposed method can reflect different audio sources.

Moreover, even though we do not have instance segmentation masks, we can provide instance-level segmentation when multiple sound sources are given.

3 PROPOSED METHOD

In the audio-visual segmentation (AVS) task, our target is to segment the corresponding sounding objects in one visual frame v given an audio a . To solve this challenging problem, we propose an instance-aware audio-visual segmentation framework by exploring cross-modal mutual correlations between visual frames and audio signals. As illustrated in Figure 2, our framework is composed of two stages, *i.e.*, a potential sounding object segmentation module and an audio-visual semantic correlation (AVSC) module.

3.1 Potential Sounding Object Segmentation

We propose the potential sounding object segmentation module to enumerate the potential sounding instances in video frames. However, due to the differences between the typical segmentation datasets and the audio-visual segmentation dataset, we cannot directly utilize the semantic segmentation models or panoptic segmentation models to achieve this goal. The difficulty is mainly induced by the specificity of the audio-visual segmentation dataset. First, different from the existing typical segmentation dataset in which each pixel has a category label, *e.g.*, COCO [15, 17] and ADE20K [44], the ground-truth labels in the audio-visual dataset are just binary masks (*i.e.*, 1 represents sounding and 0 represents silent) even though a video frame involves multiple sounding objects. Second, there are nearly 89% of the video frames that only contain one sounding object. In other words, most of the ground truths only contain one segmentation object, while the COCO or

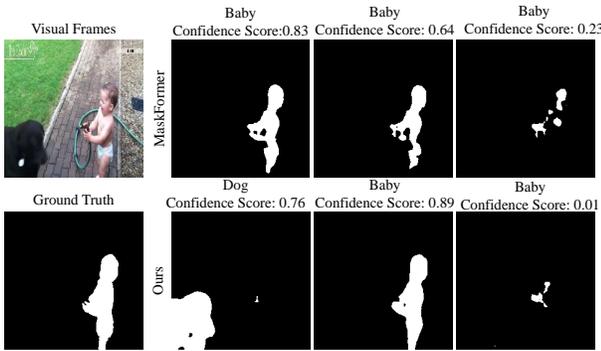


Figure 3: Visualization of instance masks obtained by MaskFormer and our segmentation network. Though we query 100 instance masks for both MaskFormer and our segmentation network, MaskFormer tends to generate masks only on one object while ours generates queries at various regions thanks to our silent object-aware segmentation loss.

ADE20K usually have multiple instances in an image. This would render a segmentation network to become a saliency segmentation framework easily, failing to propose various instances in a frame.

To solve the above problems, we employ a potential sounding object segmentation network, and the architecture design inherits the spirits of the MaskFormer series [7, 8]. As shown in Figure 2, our network is mainly composed of three components: an image encoder, a pixel decoder, and a transformer decoder. The learnable queries in the transformer decoder represent the embedding features of instances and will be mapped to binary masks and their correspondence classification scores. Note that even though we adopt the architecture of MaskFormer [7], directly training the network with the ground-truth will lead to inferior segmentation results, as illustrated in Figure 3. Therefore, we introduce our salient object-aware segmentation objective to train this network in Section 3.4, and then obtain the set of potential sounding objects even without pixel-wise label annotations.

In detail, we feed $v \in \mathbb{R}^{3 \times H \times W}$ into the segmentation network, where H and W refer to the height and width of the visual frame, respectively. After the image encoder, the image feature f is sent into the transformer decoder and the pixel decoder, respectively. The pixel decoder decodes f into an image pixel feature. The transformer decoder outputs the embedding vectors of query instances, and the embedding vectors are further mapped into the final predicted results $q = \{(p_i, m_i)\}_{i=1}^N$, where $p_i \in \Delta^{K+1}$ denotes classification confidence scores of masks, and $m_i \in \{0, 1\}^{H \times W}$ represents the binary masks. Here, K is the number of categories in the dataset. Besides, we also set a non-object category to represent the background regions. Hence the total number of categories is $K + 1$.

To mine the potential sounding objects as much as possible, the number of the predicted masks N (i.e., queries) should be larger than that of the ground-truth N_{gt} . Here, we formulate the ground-truth segments as $z_{gt} = \{(c_j, m_j) | c_j \in \{1, \dots, K\}, m_j \in \{0, 1\}^{H \times W}\}_{j=1}^{N_{gt}}$, where m_j and c_j represent a binary mask and the category label of a sounding object, respectively. During training, we utilize the bipartite matching [8] to obtain the matching index σ between the predicted results and the ground-truth, and then calculate the

training loss. After training, the potential sounding object segmentation network outputs the potential instance masks and their corresponding confidence scores. In this fashion, we could obtain the *Query Instance Mask Set* which contains the masks of instances and the confidence score of each mask. They will be employed as prior guidance to achieve the audio-visual semantic correlation in the next stage.

3.2 Audio Feature Extraction

To accomplish instance-aware audio segmentation, we need to establish a robust semantic association between the audio and visual modalities. Consequently, our primary focus lies in ensuring that the extracted audio features are semantically discriminative. Here, we employ the Bidirectional Encoder Audio Transformer (BEATs) [6], an audio pre-training model that utilizes acoustic tokenizers to enhance audio semantics, to extract pertinent audio features.

As depicted in Figure 2, we first feed the Mel-Filters transformed audio signals into the pre-trained audio encoder. Subsequently, we apply multi-layer perceptron (MLP) blocks to obtain the audio semantic distribution $P_a \in \Delta^K$ over K categories. In this context, Δ^K represents the K -dimensional probability simplex. Note that our approach does not possess knowledge of the category labels of each input audio. Instead, we only know the number of audio categories within the dataset. Hence, we formulate the instance information as the category prior guidance to guarantee the semantic distribution of the audio. To be specific, we assume that the class distribution of instances is consistent with the audio distribution. Then, we regard the probability of the audio category corresponding to the instance category id as the probability of the object making a sound. In this fashion, we could obtain a combined mask that indicates the sounding probability of each potential instance, and the ground-truth mask is the supervision to suppress the probability of silent instances and encourage the probability of sounding instances. This is the distinct difference from existing sounding object localization approaches under an audio-label supervised paradigm. The implementation details will be explained in Section 3.3.

3.3 Audio Visual Mutual Semantics Alignment

To ensure audio-visual semantic alignment, we incorporate the audio distribution into the instance masks. Specifically, under the assumption that the instance semantic distribution is the same as the audio semantic distribution, we could regard the audio category probability as the sounding probability of each instance and further obtain the audio instance aware localization map. However, the queried instance mask set procured from our segmentation network might exhibit a significant level of noise. For example, numerous masks encompass the same instance but with varying quality, and some instance masks are even classified into the non-object category. Consequently, we design a category filter and a score filter to eliminate invalid segmentation. Given the predicted results $q = (p_i, m_i)_{i=1}^N$, the category filter removes segmentation results belonging to the non-object category. Subsequently, the score filter selects the instance mask with the highest mask confidence p_i within each category. Here, we would like to emphasize that the current dataset does not contain two instances of the same category producing sound simultaneously, and this case might need extra information, such as text guidance, to localize the sounding one.

Afterwards, we obtain the *Potential Instance Mask Set* $q_r = \{c_i, m_i\}_{i=1}^{N_p}$, where $c_i \in \{0, 1\}^K$ is the one-hot vector of a certain category, and $m_i \in \{0, 1\}^{H \times W}$ is the binary mask of each instance. Then, we multiply the instance masks with their corresponding audio distribution scores $p_a \in [0, 1]$ to obtain an audio semantics-aware localization map $\mathcal{S}_{asl} \in \mathcal{R}^{H \times W}$:

$$\mathcal{S}_{asl} = \sum_{j=1}^{N_p} p_{c_j} \cdot m_j. \quad (1)$$

The audio semantics-aware localization map \mathcal{S}_{asl} indicates the potential locations of sounding objects. To establish the connection between audio and visual representations, we introduce an audio-visual correspondence loss \mathcal{L}_{avc} . If the audio instance aware localization map lies in or resembles the ground-truth mask, our loss will encourage the corresponding category score of the audio to be higher. Otherwise, the loss will penalize the audio category score. In this way, we can effectively reduce the impact of silent objects from the potential sounding object set.

3.4 Training Objective

Training Objectives for Segmentation Network. For the potential sounding objects segmentation network, our objective is to acquire the potential instance masks and the confidence scores associated with each instance. Here, we first employ an objective \mathcal{L}_{mask_cls} , composed of a cross-entropy classification loss and a binary mask loss \mathcal{L}_{mask} . The binary mask loss involves a focal loss [16] and a dice loss [22]. \mathcal{L}_{mask_cls} is defined as follows:

$$\mathcal{L}_{mask_cls}(q, z_{gt}) = \sum_{j=1}^N [-\log p_{\sigma(j)}(c_{gt}) + \mathcal{L}_{mask}(m_{\sigma(j)}, m_{gt})], \quad (2)$$

$$\mathcal{L}_{mask}(m, m_{gt}) = \lambda_f \mathcal{L}_{Focal}(m, m_{gt}) + \lambda_d \mathcal{L}_{Dice}(m, m_{gt}), \quad (3)$$

where the λ_f and λ_d are hyper-parameters and set to 20 and 1, respectively.

Moreover, to adapt the segmentation network to the audio-visual segmentation task, we design a silent object-aware segmentation loss \mathcal{L}_{soas} by reducing the overlapping regions between the non-object regions and the foreground regions, expressed by:

$$\mathcal{L}_{soas}(m_n, m_{gt}) = \sum_{j=1}^{N_n} \frac{m_{n_j} \cap (\bigcup_{k=1}^{N_{gt}} m_{gt_k})}{m_{n_j} \cup (\bigcup_{k=1}^{N_{gt}} m_{gt_k})}, \quad (4)$$

where m_{n_j} and m_{gt_k} are the predicted mask belonging to the non-object set and ground truth set respectively, and N_n and N_{gt} indicates the quantity of the non-object set and ground truth set.

To sum up, the overall loss function for the potential sounding objects segmentation network is defined by:

$$\mathcal{L}_{sn} = \mathcal{L}_{mask_cls} + \lambda_{soas} \mathcal{L}_{soas}, \quad (5)$$

where λ_{soas} is empirically set to 1.

Training Objectives for AVSC. We adopt a binary cross entropy loss (BCE) \mathcal{L}_{avc} to supervise the audio-visual semantic correlation module, formulated as:

$$\mathcal{L}_{avc} = BCE(\mathcal{S}_{asl}, M_{gt}), \quad (6)$$

where M_{gt} indicates the ground truth mask of the audio-visual input pair.

Table 1: Quantitative comparisons with the state-of-the-art on the Single-source and Multi-source sub-datasets. Results of Jaccard index (\mathcal{J}) and the F-score (\mathcal{F}) are reported.

| Settings | TPAVI [46] | | Ours | |
|---------------|------------------------|------------------------|------------------------|------------------------|
| | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ |
| Single-source | 78.7 | 87.9 | 81.29 | 88.59 |
| Multi-source | 54.0 | 64.5 | 59.5 | 65.74 |

4 EXPERIMENTS

4.1 Implementation Details

Potential sounding objects segmentation network. We employ the swin-transformer [20] pretrained on COCO [15, 17] as our backbone network. The visual input is an image of size $224 \times 224 \times 3$ pixels. The potential sounding objects segmentation network is trained using the Adam optimizer with a learning rate of $1e-4$ and a batch size of 32.

Audio-Visual Semantic Correlation Module. We adopt BEATS [6] pretrained on AudioSet [11] as our audio extractor. To achieve the audio distribution, we further employ some MLP layers. The input of the audio branch is an audio signal clip with a duration of 1 second. The sampling rate of each raw waveform is 16,000. After Mel-Filters, the audio signal is converted into spectrograms with 98×128 dimensions. In this stage, we adopt Adam optimizer with a learning rate of 0.001 and a batch size of 64. More details about the model configuration and data processing are provided in the supplementary material.

4.2 Dataset and Evaluation Metrics

AVSBench Dataset. The AVSBench dataset contains 5,356 video samples, distributed across 23 distinct auditory categories [46]. Each video has a duration of 5 seconds and is uniformly segmented into five clips. Annotations for the sounding objects are provided by the binary mask of the final frame of each video clip. Moreover, AVSBench has two distinct sub-datasets according to the number of audio sources in a video: single-source and multi-source scenarios. The single-source sub-dataset is comprised of 4,932 videos, while the multi-source one contains 424 videos. We follow the data split of [46] for training and testing.

Evaluation Metrics. We employ the Jaccard index (\mathcal{J}) [9] and F-score (\mathcal{F}) as quantitative metrics for the evaluation of the predicted binary masks of acoustic objects. The Jaccard index, \mathcal{J} , is calculated as the intersection-over-union between the predicted masks and the ground-truth. The F-score, \mathcal{F} , represents the harmonic mean of precision and recall. This metric is expressed as $\mathcal{F} = \frac{(1+\beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$, where β^2 is set to 0.3, as in [46].

4.3 Quantitative Comparison

Comparison with the AVS methods. We compare our method with the state-of-the-art method TPAVI [46] in Table 1. Table 1 indicates that our method consistently outperforms TPAVI across all the metrics. This implies our AVSC module is quite effective. In particular, in the single-source scenario, our method achieves a 2.59% improvement (from 78.7% to 81.29%) over TPAVI on the metric

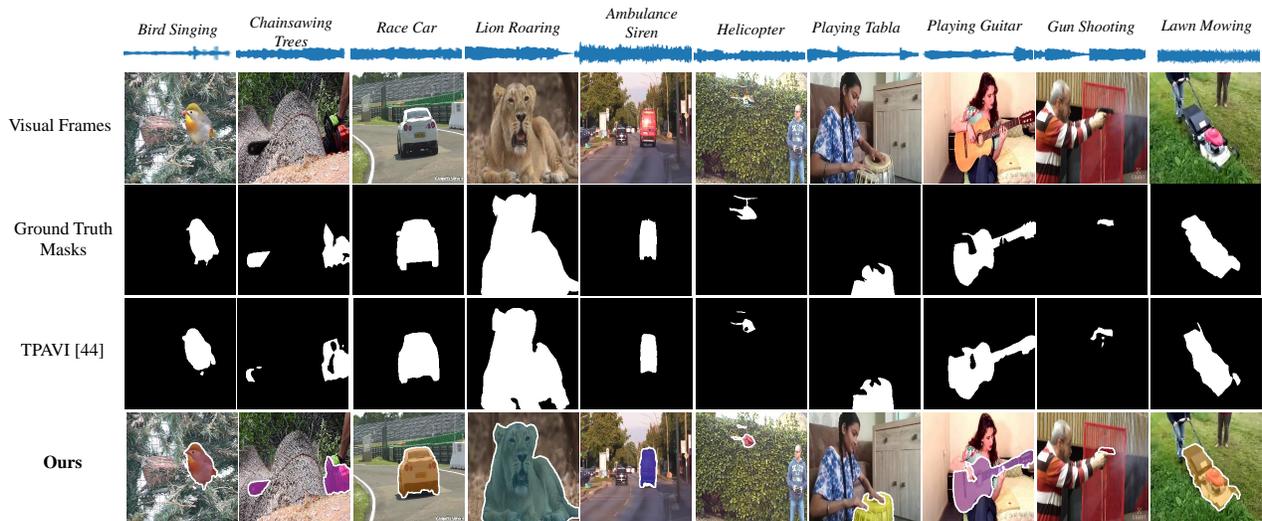


Figure 4: Qualitative comparisons with the state-of-the-art TPAVI on the single-source dataset.

Table 2: Quantitative comparisons with visual sounding localization methods on the Single-source and Multi-source sub-datasets. τ_s and τ_m represent the heatmap thresholds for optimal performance.

| Sub-dataset | EZLSL [24] | | LVS [2] | | SLAVC [23] | | SSPL [36] wo/ PCM | | SSPL [36] w/ PCM | | Ours | |
|---------------|----------------------------|---------------|----------------------------|---------------|----------------------------|---------------|----------------------------|---------------|----------------------------|---------------|---------------|---------------|
| | $(\tau_s=0.6, \tau_m=0.7)$ | | $(\tau_s=0.5, \tau_m=0.6)$ | | $(\tau_s=0.5, \tau_m=0.6)$ | | $(\tau_s=0.5, \tau_m=0.5)$ | | $(\tau_s=0.5, \tau_m=0.5)$ | | \mathcal{J} | \mathcal{F} |
| | \mathcal{J} | \mathcal{F} | | |
| Single-source | 32.99 | 44.42 | 23.28 | 30.83 | 26.53 | 35.70 | 18.72 | 32.52 | 24.04 | 36.67 | 81.29 | 88.59 |
| Multi-source | 24.01 | 27.68 | 22.73 | 31.59 | 21.40 | 23.29 | 22.83 | 25.23 | 18.67 | 23.17 | 59.5 | 65.74 |

\mathcal{J} . In the multi-Source case, our method surpasses TPAVI by a large margin of 5.5% on the metric \mathcal{J} and achieves 59.5% on the metric \mathcal{F} which is 1.24% higher than TPAVI. These results demonstrate our method performs better than the state-of-the-art on both single- and multi-source scenarios.

Comparison with the VSL methods.

We further evaluate our approach with visual sounding localization (VSL) methods, namely EZLSL [24], LVS [2], SLAVC [23], and SSPL [36]. For fair comparisons, we retrain EZLSL, SLAVC, and SSPL on the AVSBench dataset. Since LVS does not release its training code, we utilize their pre-trained model for evaluation. LVS also includes all the categories of AVSBench. Since these methods identify sounding objects with heatmaps, we further binarize their heatmaps through a series of thresholds, ranging from 0.5 to 0.9 with intervals of 0.1. Then the binary masks representing the sounding objects will be used for quantitative analysis. The results are shown in Table 2.

In Table 2, we present the optimal results for each approach under their respective thresholds. As observed, our method surpasses these VSL works by large margins. We emphasize that our exceptional performance is primarily attributed to two designs. First, we utilize a potential sounding object segmentation network to procure high-quality object masks. Second, we implement semantic association between audio and visual instances.

4.4 Qualitative Comparison

Comparison with the AVS method. Figure 4 demonstrates the results of our method and TPAVI [46] on single audio source cases. Although both methods achieve satisfactory performance, our method performs better than TPAVI near the object boundaries. For instance, as depicted in the penultimate column of Figure 4, the result of TPAVI involves a portion of the hand from the gunshot sound, while our method segments the gun more precisely.

To further illustrate the effectiveness of our audio-visual semantic correlation module, we visualize the segmentation results of the multi-source cases. As visible in Figure 7, even though the audio is a mixture of "Male Singing" and "Playing Piano", TPAVI only segments the most prominent object in the visual frames. Besides, in the multi-source setting, the segmentation masks obtained from TPAVI suffer from poor quality. This implies that TPAVI fails to capture the semantic correlation between the audio-visual pairs. In contrast, our method builds a strong semantic correlation between the audio signal and visual instances and thus achieves instance-aware sounding object segmentation.

Comparison with the VSL methods. Figure 6 illustrates the visual results of our method and the VSL methods. For instance, when the input audio signal involves "Male Singing" and "Playing Guitar", VSL methods fail to localize sounding objects in a multi-source setting. LVS and EZVSL focus on the male region while neglecting the guitar region. Although SSPL (w/o PCM) emphasizes

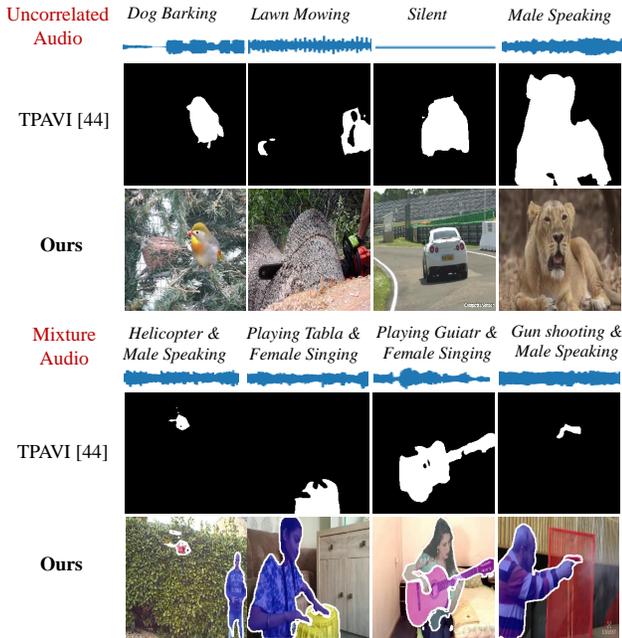


Figure 5: Visual results of changing uncorrelated and mixed audio signals for segmentation. When there is no matching instance in a video, our method does not produce any masks. The original audio-visual pairs are shown in Figure 4.

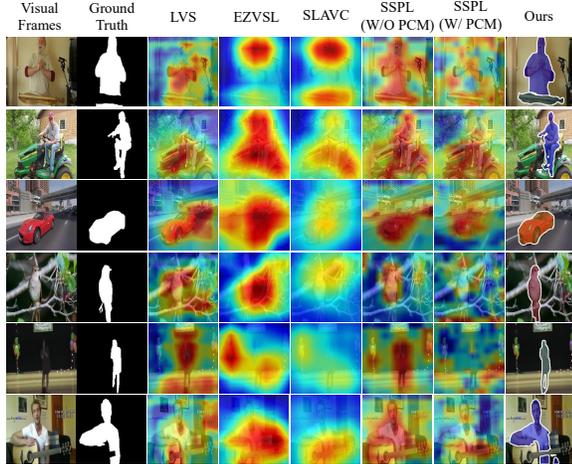


Figure 6: Visual comparison with the visual sounding localization methods including EZLSL [24], LVS [2], SLAVC [23], and SSPL [36], respectively.

the male and piano regions, it actually puts attention on all objects within the visual frames, including the microphone. Similar to SSPL (w/o PCM), SSPL (w/ PCM) also fails to localize sounding objects. In contrast, our method effectively utilizes the audio signal to accurately segment the sounding male and guitar regions, thus achieving superior performance.

Table 3: Impact of the silent object-aware segmentation loss \mathcal{L}_{soas} . \mathcal{I} represents the instance number obtained from the potential sounding objects segmentation network.

| Settings | Single-source | | | Multi-source | | |
|--------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ | $\mathcal{I} \uparrow$ | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ | $\mathcal{I} \uparrow$ |
| w/ \mathcal{L}_{soas} | 81.29 | 88.59 | 9854 | 59.5 | 65.74 | 895 |
| w/o \mathcal{L}_{soas} | 80.67 | 88.41 | 3700 | 54.35 | 63.15 | 320 |

Table 4: Impact of the AVSC on segmentation results.

| Settings | Single-source | | Multi-source | |
|----------|------------------------|------------------------|------------------------|------------------------|
| | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ | $\mathcal{J} \uparrow$ | $\mathcal{F} \uparrow$ |
| w/ AVSC | 81.29 | 88.59 | 59.5 | 65.74 |
| w/o AVSC | 77.42 | 85.62 | 53.96 | 62.05 |

4.5 Ablation Study

Impact of the silent object-aware segmentation loss. As discussed in Section 3.1, the majority of data in the audio-visual segmentation dataset only contains one sounding object, leading to ambiguity in training the segmentation network. Hence, we introduce a silent object-aware segmentation loss \mathcal{L}_{soas} to enhance the variety of segmented instances. We perform an ablation study to investigate its impact in Table 3. The results indicate that the silent object-aware segmentation loss improves audio-visual segmentation performance, particularly in the multi-source setting. Notably, the improvement of the metric \mathcal{J} is nearly 5%, and \mathcal{F} increases from 63.15% to 65.74%. To further illustrate the effectiveness of \mathcal{L}_{soas} , we analyze the segmented instance number \mathcal{I} in the test set for both single- and multi-source datasets. As shown in Table 3, the number of instances from the segmentation network with \mathcal{L}_{soas} is three times higher than that without \mathcal{L}_{soas} . This suggests that the silent object-aware segmentation loss effectively encourages our segmentation network to segment a more diverse range of instances.

Impact of AVSC. To demonstrate the effectiveness of our audio-visual semantics association, we predict masks for sounding objects without using audio information. Then, we select the object with the highest confidence score as the final audio-visual segmentation result. As illustrated in Table 4, the performance of our method decreases significantly in the absence of audio guidance (w/o AVSC). \mathcal{J} decreases by 3.87% for the single-source dataset and 5.54% for the multi-source dataset. This indicates that our method adequately exploits the audio signal in segmentation.

To further emphasize the sensitivity of our method to audio signals, we manually change the corresponding audio signal of the visual frames and report the visual results in Figure 5. Note that, the ground truth of these audio-visual pairs can be found in Figure 4. As depicted in the second and third rows, when the input audio signal is unrelated to the visual frames, TPAVI still segments the object in the image, whereas our method does not segment any object. When we mix audio signals of two instances in a visual frame, TPAVI only segments one sounding object. In contrast, our method not only successfully locates two sounding objects but also achieves

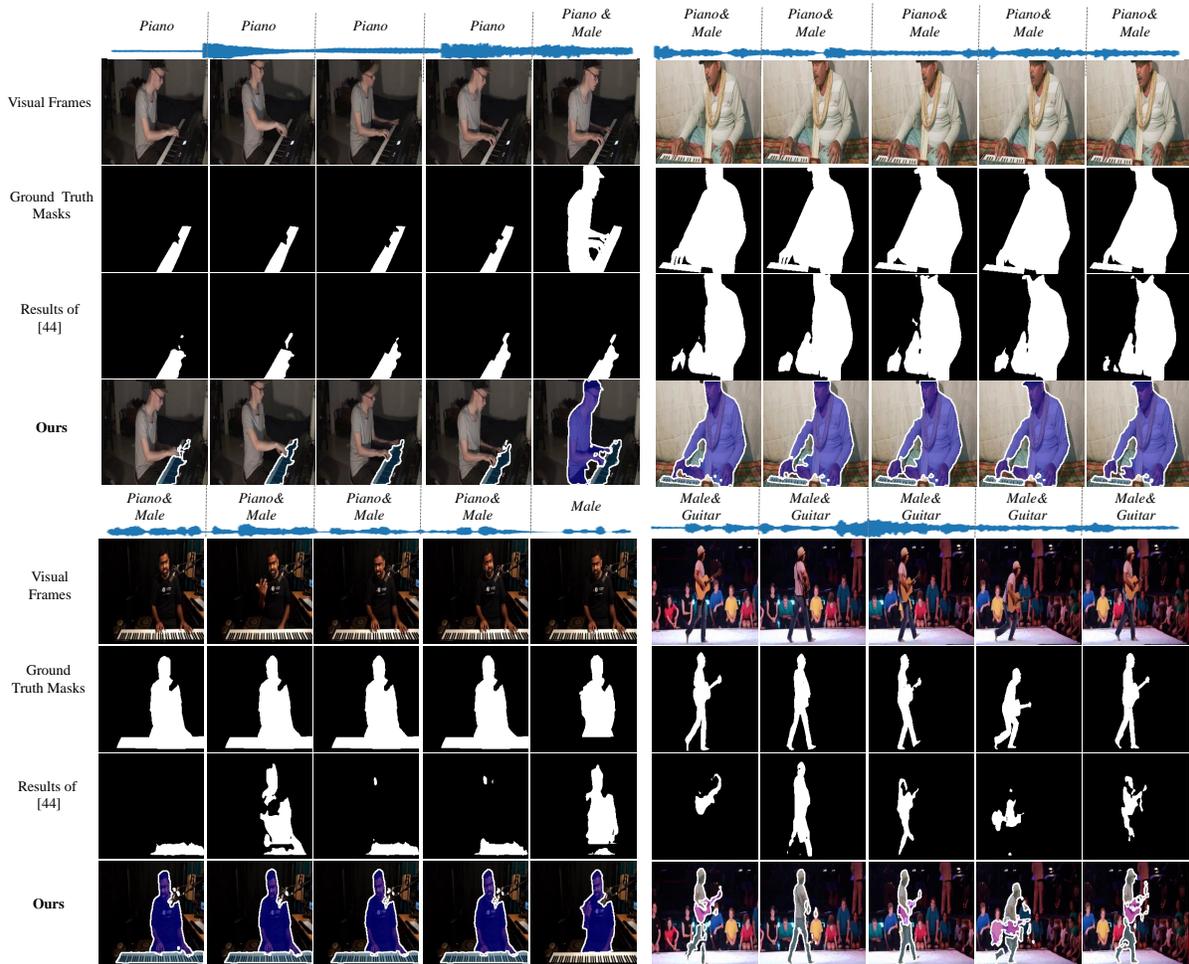


Figure 7: Qualitative comparisons with TPAVI [46] on the multi-source dataset. Our method not only successfully segments sounding objects but also provides instance-level segmentation.

sound source localization at an instance level, *i.e.*, associating sound sources with the corresponding objects.

5 CONCLUSION

In this paper, we present an audio-visual instance-aware segmentation approach that can effectively segment sounding objects according to audio signals. With the help of our silent object-aware segmentation training loss, we are able to segment all potential sounding objects in a video. After obtaining potential sounding objects, we can match object masks to audio signals. In other words, audio signals can modulate the segmentation results in an effective manner. Moreover, our audio-visual semantic association is instance-aware. Even though the segmentation ground-truth provided by the AVS benchmark do not provide instance-level annotation, our method can produce instance-level segmentation associated with multiple sound sources. Experiments on the popular used AVS benchmark demonstrate the superiority of our method over the state-of-the-art. More importantly, our method achieves the effect that segmentation results can be effectively adjusted by different audio signals.

REFERENCES

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *ECCV*. Springer, 213–229.
- [2] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. 2021. Localizing visual sounds the hard way. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16867–16876.
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* 40, 4 (2017), 834–848.
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* (2017).
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on computer vision (ECCV)*. 801–818.
- [6] Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. 2022. BEATs: Audio Pre-Training with Acoustic Tokenizers. *arXiv preprint arXiv:2212.09058* (2022).
- [7] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. 2022. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1290–1299.

- [8] Bowen Cheng, Alex Schwing, and Alexander Kirillov. 2021. Per-pixel classification is not all you need for semantic segmentation. *Advances in Neural Information Processing Systems* 34 (2021), 17864–17875.
- [9] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010. The pascal visual object classes (voc) challenge. *International journal of computer vision* 88 (2010), 303–338.
- [10] Rui Fan, Hengli Wang, Peide Cai, and Ming Liu. 2020. Sne-roadseg: Incorporating surface normal information into semantic segmentation for accurate freespace detection. In *ECCV*. Springer, 340–356.
- [11] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio Set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP*.
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [13] Dominic Henze, Klaidi Gorishti, Bernd Bruegge, and Jan-Philipp Simen. 2019. Audioforesight: A process model for audio predictive maintenance in industrial environments. In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*. IEEE, 352–357.
- [14] Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou. 2020. Discriminative sounding objects localization via self-supervised audiovisual matching. *Advances in Neural Information Processing Systems* 33 (2020), 10077–10087.
- [15] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. 2019. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9404–9413.
- [16] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.
- [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*. 740–755.
- [18] Jinxiang Liu, Chen Ju, Weidi Xie, and Ya Zhang. 2022. Exploiting Transformation Invariance and Equivariance for Self-supervised Sound Localisation. In *Proceedings of the 30th ACM International Conference on Multimedia*. 3742–3753.
- [19] Xian Liu, Rui Qian, Hang Zhou, Di Hu, Weiyao Lin, Ziwei Liu, Bolei Zhou, and Xiaowei Zhou. 2022. Visual sound localization in the wild by cross-modal interference erasing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 1801–1809.
- [20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10012–10022.
- [21] Timo Lüddecke and Alexander Ecker. 2022. Image segmentation using text and image prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7086–7096.
- [22] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*. Ieee, 565–571.
- [23] Shentong Mo and Pedro Morgado. 2022. A Closer Look at Weakly-Supervised Audio-Visual Source Localization. *arXiv preprint arXiv:2209.09634* (2022).
- [24] Shentong Mo and Pedro Morgado. 2022. Localizing visual sounds the easy way. In *ECCV*. Springer, 218–234.
- [25] Takashi Oya, Shohei Iwase, Ryota Natsume, Takahiro Itazuri, Shugo Yamaguchi, and Shigeo Morishima. 2020. Do we need sound for sound source localization?. In *Proceedings of the Asian Conference on Computer Vision*.
- [26] Xingqun Qi, Chen Liu, Lincheng Li, Jie Hou, Haoran Xin, and Xin Yu. 2023. EmotionGesture: Audio-Driven Diverse Emotional Co-Speech 3D Gesture Generation. *arXiv:2305.18891* [cs.CV]
- [27] Xingqun Qi, Chen Liu, Muyi Sun, Lincheng Li, Changjie Fan, and Xin Yu. 2023. Diverse 3D Hand Gesture Prediction from Body Dynamics by Bilateral Hand Disentanglement. *arXiv preprint arXiv:2303.01765* (2023).
- [28] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. 2020. Multiple sound sources localization from coarse to fine. In *ECCV*. 292–308.
- [29] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III* 18. Springer, 234–241.
- [30] Sebastian Schneider and Paul Wilhelm Dierkes. 2021. Localize Animal Sound Events Reliably (LASER): A new software for sound localization in zoos. *Journal of Zoological and Botanical Gardens* 2, 2 (2021), 146–163.
- [31] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. 2018. Learning to localize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4358–4366.
- [32] Arda Senocak, Hyeonngon Ryu, Junsik Kim, and In So Kweon. 2022. Learning sound localization better from semantically similar samples. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. 4863–4867.
- [33] Arda Senocak, Hyeonngon Ryu, Junsik Kim, and In So Kweon. 2022. Less can be more: Sound source localization with a classification model. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3308–3317.
- [34] Jiayin Shi and Chao Ma. 2022. Unsupervised Sounding Object Localization with Bottom-Up and Top-Down Attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1737–1746.
- [35] Ma Shuo, Yanli Ji, Xing Xu, and Xiaofeng Zhu. 2021. Vision-guided music source separation via a fine-grained cycle-separation network. In *Proceedings of the 29th ACM International Conference on Multimedia*. 4202–4210.
- [36] Zengjie Song, Yuxi Wang, Junsong Fan, Tieniu Tan, and Zhaoxiang Zhang. 2022. Self-supervised predictive learning: A negative-free method for sound source localization in visual scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3222–3231.
- [37] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. 2021. Segformer: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 7262–7272.
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [39] Jinghua Wang, Zhenhua Wang, Dacheng Tao, Simon See, and Gang Wang. 2016. Learning common and specific features for RGB-D semantic segmentation with deconvolutional networks. In *ECCV*. 664–679.
- [40] Abdelrahman Younes, Daniel Honerkamp, Tim Welschhold, and Abhinav Valada. 2023. Catch me if you hear me: Audio-visual navigation in complex unmapped environments with moving sounds. *IEEE Robotics and Automation Letters* (2023).
- [41] Tianhao Zhang, Waqas Aftab, Lyudmila Mihaylova, Christian Langran-Wheeler, Samuel Rigby, David Fletcher, Steve Maddock, and Garry Bosworth. 2022. Recent advances in video analytics for rail network surveillance for security, trespass and suicide prevention—A survey. *Sensors* 22, 12 (2022), 4324.
- [42] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. 2018. The sound of pixels. In *Proceedings of the European conference on computer vision (ECCV)*. 570–586.
- [43] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. 2017. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2881–2890.
- [44] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. 2017. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 633–641.
- [45] Jinxing Zhou, Dan Guo, and Meng Wang. 2022. Contrastive positive sample propagation along the audio-visual event line. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [46] Jinxing Zhou, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birchfield, Dan Guo, Lingpeng Kong, Meng Wang, and Yiran Zhong. 2022. Audio-Visual Segmentation. In *ECCV*. Springer, 386–403.
- [47] Xinchu Zhou, Dongzhan Zhou, Di Hu, Hang Zhou, and Wanli Ouyang. 2023. Exploiting Visual Context Semantics for Sound Source Localization. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 5199–5208.
- [48] Xinchu Zhou, Dongzhan Zhou, Wanli Ouyang, Hang Zhou, and Di Hu. 2023. SeCo: Separating Unknown Musical Visual Sounds with Consistency Guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 5168–5177.

6 APPENDICES

Table 5: Comparisons with TPAVI on the inference time (Time), computation cost (GFLOPs), and parameter amount (Param).

| Metric | Time (s) | GFLOPs | Param (MB) |
|----------------|----------|--------|------------|
| TPAVI-ResNet50 | 0.0033 | 33.878 | 91.40 |
| TPAVI-PVTv2 | 0.03994 | 30.84 | 101.32 |
| Ours-ResNet50 | 0.013395 | 13.28 | 45.52 |
| Ours-PVTv2 | 0.01724 | 29.83 | 92.94 |
| Ours-Swin-B | 0.01853 | 30.75 | 101.80 |

6.1 The details about matching index σ

The matching index σ actually represents a list that contains the index pair of the ground truth and its matching prediction. In other words, the list indicates which prediction is the most possible matching one with the ground truth. In the training stage, the matching pairs are utilized to calculate the loss, thus influencing the training effectiveness and the final results of the segmentation model. To be specific, incorrect matching pair causes inaccurate loss to update model parameters, degrading the model performance. Conversely, accurate matching facilitates the learning process of our model, thus improving the model performance.

6.2 Comparisons with AVS method

We compare the time complexity of models. As suggested by Table 5, our framework shows superiority on most of the metrics, *i.e.* GFLOPs and Parameter Amount. Note that, we also add the time for audio feature extraction into the inference time. Since TPAVI only provides the extracted audio features, we cannot take this time into account. Although our method takes longer inference time than TPAVI, it still achieves real-time interactions.

6.3 Ablation study about various backbones

To demonstrate the effectiveness of our proposed module, we utilize the same backbone models as TPAVI in our framework. More specifically, we employ ResNet50 and PVT-v2 trained on ImageNet as the backbone of our segmentation model. Concurrently, VG-Gish is adopted as the audio encoder for our implementation. As suggested by Table 7, our proposed method consistently achieves superior results on both metrics under the same backbone setting, indicating the effectiveness of our approach. While the improved backbone undoubtedly provides a solid foundation for the framework performance, the distinct advantage of our module should not be underemphasized. These modules are uniquely devised to tackle the challenge of ineffective audio and to facilitate instance-level audio-visual association, as delineated in Figure 5.

6.4 Ablation study on batchsize

In the training stage, we utilize the model with the best performance on the validation set for the next stage of training. We re-train

Table 6: Comparisons with TPAVI under the silent audio and the unmatching audio cases. In this table, RA indicates the recognition accuracy.

| Setting | | Single-source | | Multi-source | |
|------------|-------------|---------------|-------------|--------------|-------------|
| | | mIoU | RA | mIoU | RA |
| Silent | TPAVI | 81.47 | 0.0 | 88.34 | 0.0 |
| | Ours | 99.00 | 94.4 | 97.79 | 80.8 |
| Unmatching | TPAVI | 81.30 | 0.0 | 89.22 | 0.0 |
| | Ours | 93.95 | 66.9 | 96.15 | 60.9 |

our framework with the consistent batch size (20) of TPAVI. As indicated by Table 8, under the same batch size, our method is still superior to TPAVI [46]. For instance, on the multi-source sub-dataset, our method (ResNet50) achieves 1.7% Jaccard Index and 3.71% F-Score improvement over TPAVI (ResNet50). This proves that the performance improvement is mainly due to our method design rather than the batch size setting.

Table 4 represents the ablation study for the "Audio-Visual Semantic Correlation Module" (AVSC). Note that, our proposed method aims at allowing segmentation results to be controlled by different audio signals rather than overfitting to saliency objects. To demonstrate that our method solves the problem, for the *w/o* AVSC, we select the segmented instances with the highest confidence score as the predicted results. Since 89% of data in AVSBench only contain one sounding source, and the sound object is the most significant one in visual frames, Jaccard Index on both single- and multi-source datasets only increases by 3.88% and 5.54% separately.

6.5 Further analysis for AVSC

To further demonstrate the superiority of AVSC, we conduct experiments by changing the original audio signal of audio-visual pairs into the unmatching one or silent one. Under the two cases, AVS models should not segment any regions in visual frames. Note that, different from introducing Jaccard Index and F-score to measure the quality of segmented sounding regions, we focus on whether the models still segment regions in images under the two cases, *i.e.* suffer from the overfitting phenomenon. Hence, we introduce the mIoU to measure the quality of segmented silent regions, and the recognition accuracy of silent frames or audio-visual unmatching frames (The predicted masks without any segmented sounding regions are regarded as the correct recognition.).

As suggested by Table 6, our method achieves higher recognition accuracy of silent or unmatching audio-visual pairs while TPAVI tends to segment regions in the visual frames even for the silent cases. This further demonstrates our proposed method is sensitive to audio changes, solving the overfitting problem.

Table 7: Effect of the various backbones. For the segmentation network, we employ ResNet50 and PVT-v2 trained on ImageNet as the backbone. For the audio branch, VGGish is adopted as the audio encoder for our implementation

| Metric | Setting | TPAVI [46] | | Ours | |
|----------|---------------|-----------------|---------------|----------------------|----------------------|
| | | ResNet50+VGGish | PVT-v2+VGGish | ResNet50+VGGish | PVT-v2+VGGish |
| <i>J</i> | Single-Source | 72.79 | 78.74 | 77.02 (+4.23) | 80.57 (+1.83) |
| | Multi-Source | 47.88 | 54.00 | 49.58 (+1.7) | 58.22 (+4.22) |
| <i>F</i> | Single-Source | 84.80 | 87.90 | 85.24 (+0.44) | 88.19 (+0.29) |
| | Multi-Source | 57.80 | 64.50 | 61.51 (+3.71) | 65.10 (+0.6) |

Table 8: Effects about the batchsize. In this experiment, we adopt the same batchsize as TPAVI [46] to further demonstrate the effectiveness of our method.

| Setting | Metric | TPAVI (ResNet50) | Ours (ResNet50) | TPAVI (PVT-v2) | Ours (PVT-v2) | Ours (Swin-B) |
|---------------|----------|------------------|----------------------|----------------|----------------------|---------------|
| Single-Source | <i>J</i> | 72.79 | 77.02 (+4.23) | 78.74 | 80.57 (+1.83) | 81.12 |
| | <i>F</i> | 84.80 | 85.24 (+0.44) | 87.90 | 88.19 (+0.29) | 88.37 |
| Multi-Source | <i>J</i> | 47.88 | 49.58 (+1.7) | 54.00 | 58.22 (+4.22) | 59.04 |
| | <i>F</i> | 57.8 | 61.51 (+3.71) | 64.5 | 65.10 (+0.6) | 64.02 |