

Cyclic Co-Learning of Sounding Object Visual Grounding and Sound Separation

Yapeng Tian¹, Di Hu^{2,3*}, Chenliang Xu^{1*}

¹University of Rochester, ²Gaoling School of Artificial Intelligence, Renmin University of China

³ Beijing Key Laboratory of Big Data Management and Analysis Methods

{yapengtian, chenliang.xu}@rochester.edu, dihu@ruc.edu.cn

Abstract

There are rich synchronized audio and visual events in our daily life. Inside the events, audio scenes are associated with the corresponding visual objects; meanwhile, sounding objects can indicate and help to separate their individual sounds in the audio track. Based on this observation, in this paper, we propose a cyclic co-learning (CCoL) paradigm that can jointly learn sounding object visual grounding and audio-visual sound separation in a unified framework. Concretely, we can leverage grounded object-sound relations to improve the results of sound separation. Meanwhile, benefiting from discriminative information from separated sounds, we improve training example sampling for sounding object grounding, which builds a co-learning cycle for the two tasks and makes them mutually beneficial. Extensive experiments show that the proposed framework outperforms the compared recent approaches on both tasks, and they can benefit from each other with our cyclic co-learning. The source code and pre-trained models are released in <https://github.com/YapengTian/CCOL-CVPR21>.

1. Introduction

Seeing and hearing are two of the most important senses for human perception. Even though the auditory and visual information may be discrepant, the percept is unified with multisensory integration [5]. Such phenomenon is considered to be derived from the characteristics of specific neural cell, as the researchers in cognitive neuroscience found the superior temporal sulcus in the temporal cortex of the brain can simultaneously response to visual, auditory, and tactile signal [18, 40]. Accordingly, we tend to perform as unconsciously correlating different sounds and their visual producers, even in a noisy environment. For example, for a cocktail-party scenario contains multiple sounding and silent instruments as shown in Fig. 1, we can effortlessly filter out the silent ones and identify different sounding ob-

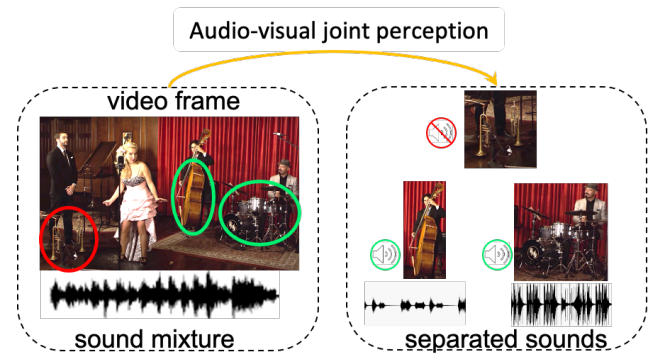


Figure 1. Our model can perform audio-visual joint perception to simultaneously identify silent and sounding objects and separate sounds for individual audible objects.

jects, and simultaneously separate the sound for each playing instrument, even faced with a *static visual image*.

For computational models, the multi-modal sound separation and sounding object alignment capacities reflect in audio-visual sound separation (AVSS) and sound source visual grounding (SSVG), respectively. AVSS aims to separate sounds for individual sound sources with help from visual information, and SSVG tries to identify objects that make sounds in visual scenes. These two tasks are primarily explored isolatedly in the literature. Such a disparity to human perception motivates us to address them in a co-learning manner, where we leverage the joint modeling of the two tasks to discover objects that make sounds and separate their corresponding sounds without using annotations.

Although existing works on AVSS [8, 32, 11, 49, 13, 46, 9] and SSVG [27, 38, 43, 3, 33] are abundant, it is non-trivial to jointly learn the two tasks. Previous AVSS methods implicitly assume that all objects in video frames make sounds. They learn to directly separate sounds guided by encoded features from either entire video frames [8, 32, 49] or detected objects [13] without parsing which are sounding or not in unlabelled videos. At the same time, the SSVG methods mostly focus on the simple scenario with single sound source, barely exploring the realistic cocktail-party environment [38, 3]. Therefore, these methods blindly use information from silent objects to guide separation learn-

*Corresponding authors.

ing, also blindly use information from sound separation to identify sounding objects.

Toward addressing the drawbacks and enabling the co-learning of both tasks, we introduce a new sounding object-aware sound separation strategy. It targets to separate sounds guided by only sounding objects, where the audio-visual scenario usually consists of multiple sounding and silent objects. To address this challenging task, the SSVG can jump in to help identify each sounding object from a mixture of visual objects, whose objective is unlike previous approaches that make great efforts on improving the localization precision of sound source in simple audiovisual scenario [38, 43, 3, 32]. Accordingly, it is challenging to discriminatively discover isolated sounding objects inside scenarios via the predicted audible regions visualized by heatmaps [38, 3, 20]. For example, two nearby sounding objects might be grouped together in a heatmap and we have no good principle to extract individual objects from a single located region.

To enable the co-learning, we propose to directly discover individual sounding objects in visual scenes from visual object candidates. With the help of grounded objects, we learn sounding object-aware sound separation. Clearly, a good grounding model can help to mitigate learning noise from silent objects and improve separation. However, causal relationship between the two tasks cannot ensure separation can further enhance grounding, because they only loosely interacted during sounding object selection. To alleviate the problem, we use separation results to help sample more reliable training examples for grounding. It makes the co-learning in a cycle and both grounding and separation performance will be improved, as illustrated in Fig. 2. Experimental results show that the two tasks can be mutually beneficial with the proposed cyclic co-learning, which leads to noticeable performance and outperforms recent methods on sounding object visual grounding and audio-visual sound separation tasks.

The main contributions of this paper are as follows: (1) We propose to perform sounding object-aware sound separation with the help of visual grounding task, which essentially analyzes the natural but previously ignored cocktail-party audiovisual scenario. (2) We propose a cyclic co-learning framework between AVSS and SSVG to make these two tasks mutually beneficial. (3) Extensive experiments and ablation study validate that our models can outperform recent approaches, and the tasks can benefit from each other with our cyclic co-learning.

2. Related Work

Sound Source Visual Grounding: Sound source visual grounding aims to identify the visual objects associating to specific sounds in the daily audiovisual scenario. This task is closely related to the visual localization problem of

sound source, which targets to find pixels that are associated with specific sounds. Early works in this field use canonical correlation analysis [27] or mutual information [17, 19] to detect visual locations that make the sound in terms of localization or segmentation. Recently, deep audio-visual models are developed to locate sounding pixels based on audio-visual embedding similarities [3, 32, 20, 22], cross-modal attention mechanisms [38, 43, 1], vision-to-sound knowledge transfer [10], and sounding class activation mapping [33]. These learning fashions are capable of predicting audible regions by showing heat-map visualization of single sound source in the simple scenario, but cannot explicitly detect multiple isolated sounding objects when they are sounding at the same time. In the most recent work, Hu *et al.* [21] propose to discriminatively identify sounding and silent object in the cocktail-party environment, but relying on reliable visual knowledge of objects learned from manually selected single source videos. Unlike the previous methods, we focus on finding each individual sounding object from cocktail scenario of multiple sounding and silent objects without any human annotations, and is cooperatively learned with audio-visual sound separation task.

Audio-Visual Sound Separation: Respecting for the long-history research on sound source separation in signal processing, we only survey recent audio-visual sound source separation methods [8, 32, 11, 49, 13, 46, 48, 37, 9] in this section. These approaches separate visually indicated sounds for different audio sources (*e.g.*, speech in [8, 32], music instruments in [49, 13, 48, 9], and universal sources in [11, 37]) with a commonly used mix-and-separate strategy to build training examples [16, 24]. Concretely, they generate the corresponding sounds w.r.t. the given visual objects [49] or object motions in the video [48], where the objects are assumed to be sounding during performing separation. Hence, if a visual object belonging to the training audio source categories is silent, these models usually fail to separate a all-zero sound for it. Due to the commonly existing audio spectrogram overlapping phenomenon, these methods will introduce artifacts during training and make wrong predictions during inference. Unlike previous approaches, we propose to perform sounding object-aware audio separation to alleviate the problem. Moreover, sounding object visual grounding is explored in a unified framework with the separation task.

Audio-Visual Video Understanding: Since auditory modality containing synchronized scenes as the visual modality is widely available in videos, it attracts a lot of interests in recent years. Besides sound source visual grounding and separation, a range of audio-visual video understanding tasks including audio-visual action recognition [14, 26, 29], audio-visual event localization [43, 31, 45], audio-visual video parsing [42], cross-modal generation [6, 7, 12, 51, 50, 47], and audio-visual video caption-

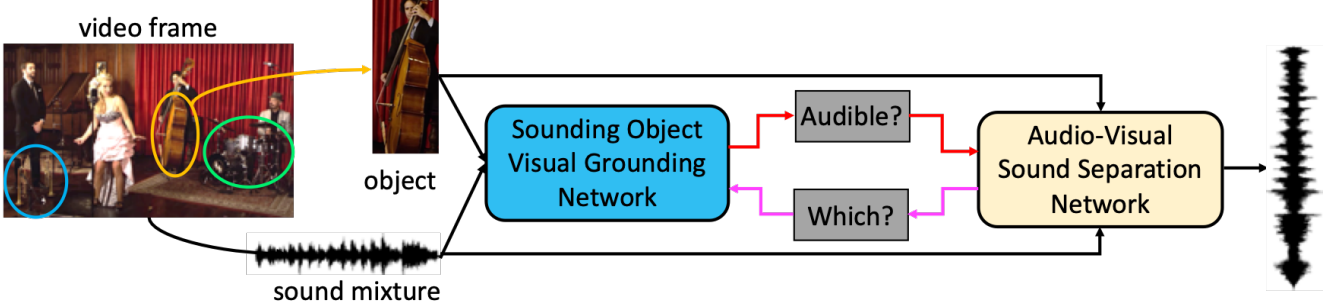


Figure 2. Cyclic Co-learning of sounding object visual grounding and audio-visual sound separation, enabled by sounding object-aware sound separation. Given detected objects from a video and the sound mixture, our model can recognize whether they are audible via a grounding network (**Audible?**) and separate their sounds with an audio-visual sound separation network to help determine potential sounding and silent sources (**Which?**).

ing [35, 41, 44] have been explored. Different from these works, we introduce a cyclic co-learning framework for both grounding and separation tasks and show that they can be mutually beneficial.

3. Method

3.1. Overview

Given an unlabeled video clip V with the synchronized sound $s(t)$ ², $\mathcal{O} = \{O_1, \dots, O_N\}$ are N detected objects in the video frames and the sound mixture is $s(t) = \sum_{n=1}^N s_n(t)$. Here, $s_n(t)$ is the separated sound of the object O_n . When it is silent, $s_n(t) = 0$. Our co-learning aims to recognize each sounding object O_n and then separate its sound $s_n(t)$ for the object.

The framework, as illustrated in Fig. 2, mainly contains two modules: sounding object visual grounding network and audio-visual sound separation network. The sounding object visual grounding network can discover isolated sounding objects from object candidates inside video frames. We learn the grounding model from sampled positive and negative audio-visual pairs. To learn sound separation in the framework, we adopt a commonly used mix-and-separate strategy [13, 16, 49] during training. Given two training video and sound pairs $\{V^{(1)}, s^{(1)}(t)\}$ and $\{V^{(2)}, s^{(2)}(t)\}$, we obtain a mixed sound:

$$s_m(t) = s^{(1)}(t) + s^{(2)}(t) = \sum_{n=1}^{N_1} s_n^{(1)}(t) + \sum_{n=1}^{N_2} s_n^{(2)}(t), \quad (1)$$

and find object candidates $\mathcal{O}^{(1)} = \{O_1^{(1)}, \dots, O_{N_1}^{(1)}\}$ and $\mathcal{O}^{(2)} = \{O_1^{(2)}, \dots, O_{N_2}^{(2)}\}$ from the two videos. The sounding object visual grounding network will recognize audible objects from $\mathcal{O}^{(1)}$ and $\mathcal{O}^{(2)}$ and the audio-visual sound separation network will separate sounds for the grounded objects. Sounds are processed in a Time-Frequency space with the short-time Fourier transform (STFT).

² $s(t)$ is a time-discrete audio signal.

With the sounding object-aware sound separation, we can co-learn the two tasks and improve audio-visual sound separation with the help of sounding object visual grounding. However, the separation performance will highly rely on the grounding model and the grounding task might not benefit from co-learning training due to the weak feedback from separation. To simultaneously evolve the both models, we propose a cyclic co-learning strategy as illustrated in Fig. 3, which has an additional backward process that utilizes separation results to directly improve training sample mining for sounding object visual grounding. In this manner, we can make the two tasks mutually beneficial.

3.2. Sounding Object Visual Grounding

Videos contain various sounds and visual objects, and not all objects are audible. To find sounding objects in videos $V^{(1)}$ and $V^{(2)}$ and further utilize grounding results for separation, we formulate sounding object visual grounding as a binary matching problem.

Sounding Object Candidates: To better support the audio-visual matching problem, we choose to follow the widely-adopted image representation strategy of visual object proposal in image captioning [25, 2], which has been also employed in the recent work on audio-visual learning [13]. Concretely, the potential audible visual objects are first proposed from videos using an object detector. In our implementation, we use the Faster R-CNN [36] object detector trained on Open Images dataset [30] from [13] to detect objects from video frames in $V^{(1)}$ and $V^{(2)}$ and obtain $\mathcal{O}^{(1)} = \{O_1^{(1)}, \dots, O_{N_1}^{(1)}\}$ and $\mathcal{O}^{(2)} = \{O_1^{(2)}, \dots, O_{N_2}^{(2)}\}$. Next, we learn to recognize sounding objects in $\mathcal{O}^{(1)}$ and $\mathcal{O}^{(2)}$ associated with $s^{(1)}(t)$ and $s^{(2)}(t)$, respectively. For simplicity, we use an object O and a sound $s(t)$ as an example to illustrate our grounding network.

Audio Network: Raw waveform $s(t)$ is transformed to an audio spectrogram S with the STFT. An VGG [39]-like 2D CNN network: VGGish followed by a global max pooling (GMP) is used to extract an audio embedding f_s from S .

Visual Network: The visual network extracts features from

detected visual object O . We use the pre-trained ResNet-18 [15] model before the last fully-connected layer to extract a visual feature map and perform a GMP to obtain a visual feature vector f_o for O .

Grounding Module: The audio-visual grounding module takes audio feature f_s and visual object feature f_o as inputs to predict whether the visual object O is one of the sounding makers for $s(t)$. We solve it using a two-class classification network. It first concatenates f_s and f_o and then uses a 3-layer Multi-Layer Perceptron (MLP) with a Softmax function to output a probability score $g(s(t), O) \in \mathcal{R}^2$. Here, if $g(s(t), O)[0] > 0.5$, $\{s(t), O\}$ is a positive pair and $s(t)$ and O are matched; otherwise, O is not a sound source.

Training and Inference: To train the sounding object visual grounding network, we need to sample positive/matched and negative/mismatched audio and visual object pairs. It is straightforward to obtain negative pairs with composing audio and objects from different videos. For example, $s^{(1)}(t)$ from $V^{(1)}$ and an randomly selected object $O_r^{(2)}$ from $V^{(2)}$ can serve as a negative pair. However, positive audio-visual pairs are hard to extract since not all objects are audible in videos. If an object from $V^{(1)}$ is not audio source, the object and $s^{(1)}(t)$ will be a negative pair, even though they are from the same video. To address the problem, we cast the positive sample mining as a multiple instance learning problem and sample the most confident pair as a positive sample with a grounding loss as the measurement:

$$\hat{n} = \operatorname{argmin}_n f(g(s^{(1)}(t), O_n^{(1)}), y_{pos}), \quad (2)$$

where $f(\cdot)$ is a cross-entropy function; $y_{pos} = [1, 0]$ is an one-hot encoding for positive pairs; $O_{\hat{n}}^{(1)}$ and $s^{(1)}$ will be the positive audio-visual pair for training. With the sampled negative and positive data, we can define the loss function to learn the sounding object visual grounding:

$$l_{grd_s} = f(g(s^{(1)}(t), O_r^{(2)}), y_{neg}) + f(g(s^{(1)}(t), O_{\hat{n}}^{(1)}), y_{pos}), \quad (3)$$

where $y_{neg} = [0, 1]$ is the negative label. The visual grounding network can be end-to-end optimized with sampled training pairs via l_{grd_s} .

During inference, we can feed audio-visual pairs $\{O_i^{(1)}, s^{(1)}(t)\}_{i=1}^{N_1}$ and $\{O_i^{(2)}, s^{(2)}(t)\}_{i=1}^{N_2}$ into the trained model to find sounding objects inside the two videos. To facilitate audio-visual sound separation, we need to detect sounding objects from the sound mixture $s_m(t)$ rather than $s^{(1)}(t)$ and $s^{(2)}(t)$, since the individual sounds are unavailable at a testing stage for separation task.

3.3. Sounding Object-Aware Separation

Given detected objects in $\mathcal{O}^{(1)}$ and $\mathcal{O}^{(2)}$, we separate sounds for each object from the sound mixture $s_m(t)$ and mute separated sounds of silent objects.

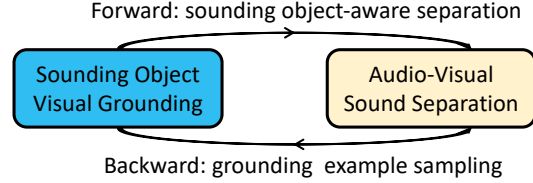


Figure 3. Cyclic co-learning. Facilitated by sounding object visual grounding, our model can employ sounding object-aware sound separation to improve separation. Meanwhile, separation results can help to do effective training sample mining for grounding.

Using an audio-visual sound separation network, we can predict sound spectrograms $\{S_n^{(1)}\}_{n=1}^{N_1}$ and $\{S_n^{(2)}\}_{n=1}^{N_2}$ for objects in $\mathcal{O}^{(1)}$ and $\mathcal{O}^{(2)}$, respectively. According to waveform relationship in Eq. 1, we can approximate spectrogram magnitude relationship as: $S^{(1)} \approx \sum_{n=1}^{N_1} S_n^{(1)}$ and $S^{(2)} \approx \sum_{n=1}^{N_2} S_n^{(2)}$. To learn the separation network, we can optimize it with a L1 loss function:

$$l_{sep} = \|S^{(1)} - \sum_{n=1}^{N_1} S_n^{(1)}\|_1 + \|S^{(2)} - \sum_{n=1}^{N_2} S_n^{(2)}\|_1. \quad (4)$$

However, not all objects are audible and spectrograms from different objects contain overlapping content. Therefore, even an object $O_n^{(1)}$ is not sounding, it can also separate non-zero sound spectrogram $S_n^{(1)}$ from the spectrogram of sound mixture S_m , which will introduce errors during training. To address the problem, we propose a sounding object-aware separation loss function:

$$l_{sep}^* = \|S^{(1)} - \sum_{n=1}^{N_1} g^*(s_m(t), O_n^{(1)}) S_n^{(1)}\|_1 + \|S^{(2)} - \sum_{n=1}^{N_2} g^*(s_m(t), O_n^{(2)}) S_n^{(2)}\|_1, \quad (5)$$

where $g^*(\cdot)$ is a binarized value of $g(\cdot)[0]$. If an object $O_n^{(1)}$ is not a sound source, $g^*(s_m(t), O_n^{(1)})$ will be equal to zero. Thus, the sounding object-aware separation can help to reduce training errors from silent objects in Eq. 4.

In addition, we introduce additional grounding loss terms to guide the grounding model learning from the sound mixture. Since we have no sounding object annotations, we adopt a similar positive sample mining strategy as in Eq. 2 and define a grounding loss as follows:

$$l_{grd_m} = \sum_{k=1}^2 \min_n f(g(s_m(t), O_n^{(k)}), y_{pos}). \quad (6)$$

3.4. Co-learning in a Cycle

Combing grounding and separation loss terms, we can learn the two tasks in a unified way with a co-learning objective function: $l_{col} = l_{grd_s} + l_{sep}^* + l_{grd_m}$.

Although our sounding object visual grounding and audio-visual sound separation models can be learned together, the two tasks loosely interact in l_{sep}^* . Clearly, a good grounding network can help improve the separation task. However, the grounding task might not be able to benefit from co-learning training since there is no strong feedback from separation to guide learning the grounding model. To further strengthen the interaction between the two tasks, we propose a cyclic co-learning strategy, which can make them benefit from each other.

If an object $O_n^{(k)}$ makes sound in video $V^{(k)}$, the separated spectrogram $S_n^{(k)}$ should be close to $S^{(k)}$; otherwise, the difference between $S_n^{(k)}$ and $S^{(k)}$ should be larger than a separated sound spectrogram from a sounding object and $S^{(k)}$. We use L_1 distance to measure dissimilarity of spectrograms: $d_n^{(k)} = \|S_n^{(k)} - S^{(k)}\|_1$, where $d_n^{(k)}$ will be small for a sounding object $O_n^{(k)}$. Based on the observation, we select the object $O_n^{(k)}$ with the minimum $d_n^{(k)}$ make the dominant sound in V^k to compose positive samples for sounding object visual grounding. Let $\hat{n}_1 = \operatorname{argmin}_n d_n^{(1)}$ and $\hat{n}_2 = \operatorname{argmin}_n d_n^{(2)}$. We can re-formulate grounding loss terms as:

$$l_{grd_s}^* = f(g(s^{(1)}(t), O_r^{(2)}), y_{neg}) + f(g(s^{(1)}(t), O_{\hat{n}_1}^{(1)}), y_{pos}),$$

$$l_{grd_m}^* = f(g(s_m(t), O_{\hat{n}_1}^{(1)}), y_{pos}) + f(g(s_m(t), O_{\hat{n}_2}^{(2)}), y_{pos}).$$

In addition, if $d_n^{(k)}$ is very large, the object $O_n^{(k)}$ is very likely not be audible, which can help us sample potential negative examples for mixed sound grounding. Specifically, we select the objects that are associated with the largest $d_n^{(k)}$, and $d_n^{(k)}$ must be larger than a threshold ϵ . Let $n_1^* = \operatorname{argmax}_n d_n^{(1)}$, s.t. $d_n^{(1)} > \epsilon$ and $n_2^* = \operatorname{argmax}_n d_n^{(2)}$, s.t. $d_n^{(2)} > \epsilon$. We can update $l_{grd_m}^*$ with learning from potential negative samples: $\hat{l}_{grd_m}^* = \sum_{k=1}^2 (f(g(s_m(t), O_{\hat{n}_k}^{(k)}), y_{pos}) + f(g(s_m(t), O_{n_k^*}^{(k)}), y_{neg}))$. Finally, we can co-learn the two tasks in a cycle with optimizing the joint cyclic co-learning loss function: $l_{ccol} = l_{grd_s}^* + l_{sep}^* + \hat{l}_{grd_m}^*$. Inside cyclic co-learning as illustrated in Fig. 3, we use visual grounding to improve sound separation and enhance visual grounding based on feedback from sound separation. The learning strategy can make the tasks help each other in a cycle and significantly improve performance for both tasks.

4. Experiments

4.1. Experimental Setting

Dataset: In our experiments, 520 online available musical solo videos from the widely-used MIT MUSIC dataset [49] is used. The dataset includes 11 musical instrument categories: accordion, acoustic guitar, cello, clarinet, erhu, ute, saxophone, trumpet, tuba, violin, and xylophone.

Methods	OTS [3]	DMC [20]	Grounding only	CoL	CCoL
Single Sound	58.7	65.3	72.0	67.0	84.5
Mixed Sound	51.8	52.6	61.4	58.2	75.9

Table 1. Sounding object visual grounding performance (%). Top-2 results are highlighted.

The dataset is relatively clean and sounding instruments are usually visible in videos. We split it into training/validation/testing sets, which have 468/26/26 videos from different categories, respectively. To train and test our cyclic co-learning model, we randomly select three other videos for each video to compose training and testing samples. Let’s denote the four videos as A, B, C, D. We compose A, B together as $V^{(1)}$ and C, D together as $V^{(2)}$, while sounds of $V^{(1)}$ and $V^{(2)}$ are only from A and C, respectively. Thus, objects from B and D in the composed samples are inaudible. Finally, we have 18,720/260/260 composed samples in our training/val/test sets for the two tasks.

Evaluation Metrics: For sounding object visual grounding, we feeding detected audible and silent objects in videos into different grounding models and evaluate their binary classification accuracy. We use the mir eval library [34] to measure sound separation performance in terms of two standard metrics: Signal-to-Distortion Ration (SDR) and Signal-to-Interference Ratio (SIR).

Implementation Details: We sub-sample audio signals at 11kHz, and each video sample is approximately 6 seconds. The STFT is calculated using a Hann window size of 1022 and a hop length of 256 and each 1D audio waveform is transformed to a 512×256 Time-Frequency spectrogram. Then, it is re-sampled to $T, F = 256$. The video frame rate is set as $1fps$ and we randomly select 3 frames per 6s video. Objects extracted from video frames are resized to 256×256 and then randomly cropped to 224×224 as inputs to our network. ϵ is set to 0.1. We use a soft sound spectrogram masking strategy as in [13, 49] to generate individual sounds from audio mixtures and adopt a audio-visual sound separation network from [49]. More details about the network can be found in our appendix. Our networks are optimized by Adam [28]. Since the sounding object visual grounding and audio-visual sound separation tasks are mutually related, we need to learn good initial models for making them benefit from cyclic co-learning. To this end, we learn our CCoL model with three steps in a curriculum learning [4] manner. Firstly, we train the sounding object visual grounding network with l_{grd_s} . Secondly, we co-learn the grounding network initialized with pre-trained weights and the separation network optimized by l_{col} . Thirdly, we use l_{ccol} to further fine-tune the two models.

4.2. Sounding Object Visual Grounding

We compare our methods to two recent methods: OTS [3] and DMC [20]. In addition, we make an ablation

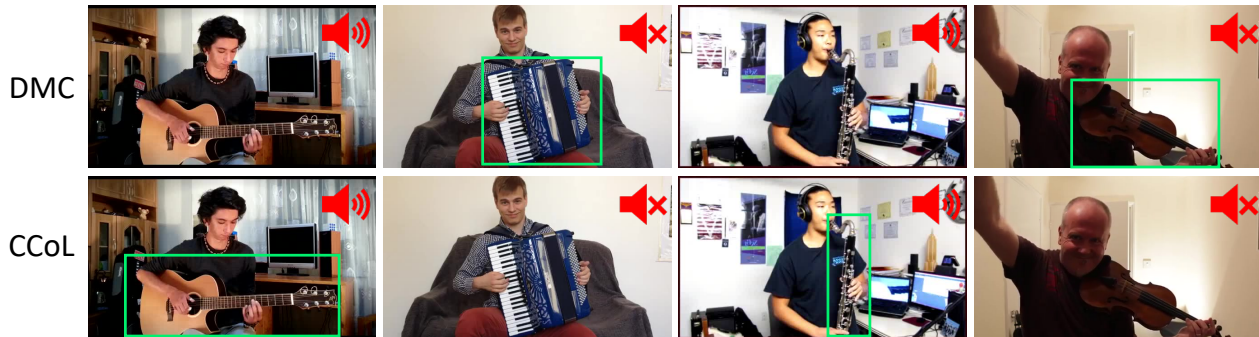


Figure 4. Qualitative results of sounding object visual grounding for both audible and silent objects. We use two icons to denote whether objects in video frames are audible or not and grounded sounding objects from DMC and CCoL are shown in green boxes. Our CCoL model can effectively identify both sounding and silent objects, while the DMC fails. Note that 2-sound mixtures are used as inputs.

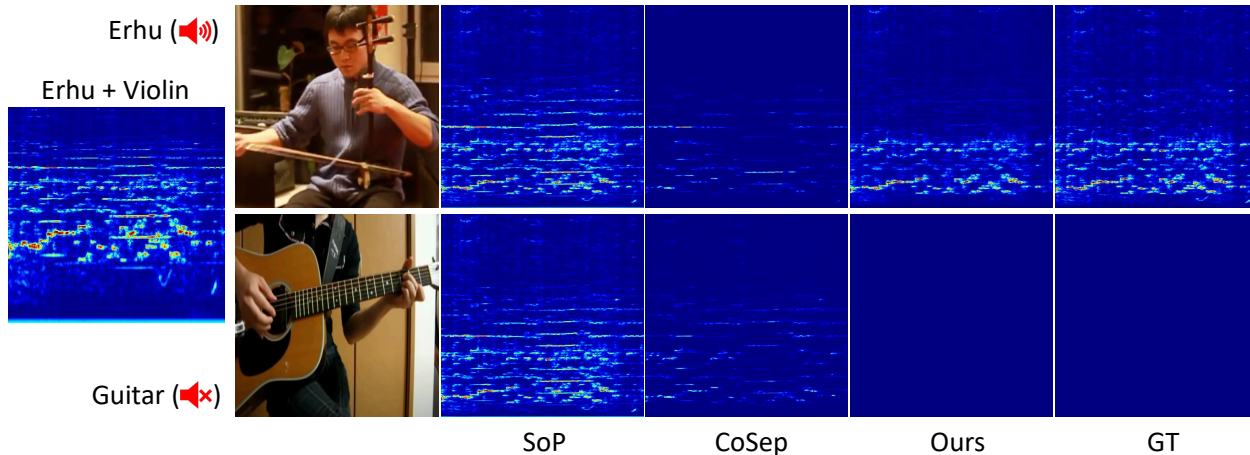


Figure 5. Qualitative results of audio-visual sound separation for both audible and silent objects. Our CCoL model can well mitigate learning noise from silent objects during training and generate more accurate sounds.

study to investigate the proposed models. The Grounding only model is trained only with grounding losses: l_{grd_s} ; the co-learning (CoL) model jointly learn visual grounding and sound separation using the l_{col} ; and the cyclic co-learning (CCoL) further strengthens the interaction between the two tasks optimized via l_{ccol} . We evaluate sounding object visual grounding performance on both solo and mixed sounds.

Table 1 and Figure 4 show quantitative and qualitative sounding object visual grounding results, respectively. Even our grounding only has already outperformed the OTS and DMC, which can validate the effectiveness of the proposed positive sample mining approach. Then, we can see that the CoL with jointly learning grounding and separation achieves worse performance than the Grounding only model. It demonstrates that the weak interaction inside CoL cannot let the grounding task benefit from the separation task. However, using separation results to help the grounding example sampling, our CCoL is significantly superior over both Grounding only and CoL models. The results can demonstrate the sounding object visual grounding can benefit from separation with our cyclic learning.

Methods	RPCA [23]	SoP [49]	CoSep [13]	Random Obj	CoL	CCoL	Oracle
SDR	-0.48	3.42	2.04	4.20	6.50	7.27	7.71
SIR	3.13	4.98	6.21	6.90	11.81	12.77	11.42

Table 2. Audio-visual sound separation performance. Top-2 results are highlighted.

4.3. Audio-Visual Sound Separation

To demonstrate the effectiveness of our CCoL framework on audio-visual sound separation, we compare it to a classical factorization-based method: RPCA [23] and two recent state-of-the-art methods: SoP [49], CoSep [13], and baselines: Random Obj and CoL in Tab. 2. The SoM [48] and Music Gesture [9] address music sound separation by incorporating dynamic visual motions, and show promising results. However, as the SoP and CoSep, they also did not consider the silent object problem. Meanwhile, since there are no source code for the two approaches, we will not include them into comparison. *Note that SoP and CoSep are trained using source code provided by the authors and the same training data (including audio preprocessing) as ours.* Moreover, we show separation results of an Oracle, which feeds ground truth grounding labels of mixed sounds to train the audio-visual separation network.

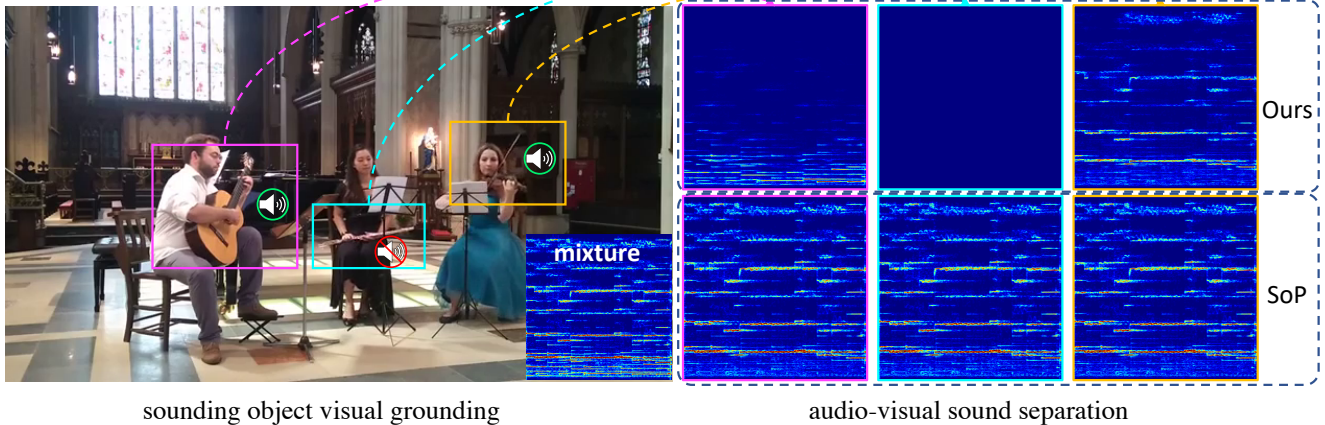


Figure 6. Real-world sounding object visual grounding and audio-visual sound separation. The *guitar* and *violin* are playing and the *flute* is visible but not audible. Our model can identify sounding objects: *guitar* and *violin* and silent object: *flute*. Moreover, it can simultaneously separate individual sounds for each instrument, while the SoP using the same noisy training data as ours fails to associate objects with the corresponding sounds, thus obtains poor separation results.

We can see that our CoL outperforms the compared SoP, CoSep, and Random Obj, and CCoL is better than CoL. The results demonstrate that sounding object visual grounding in the co-learning can help to mitigate training errors from silent video objects in separation, and separation performance can be further improved with the help of enhanced grounding model by cyclic co-learning. Compared to the Oracle model, it is reasonable to see that CCoL has slightly lower SDR. A surprising observation is that CoL and CCoL achieve better results in terms of SIR. One possible reason is that our separation networks can explore various visual objects as inputs during joint grounding and separation learning, which might make the models more robust on SIR.

Moreover, we illustrate quantitatively results for testing videos with both audible and silent objects in Tab. 3. Both SoP and CoSep are blind to whether a visual object makes sound and they will generate non-zero audio waveform for silent objects, and the Random Obj is limited in identifying object silent and sounding objects, thus they will have even lower SDRs and SIRs. However, both CoL and CCoL are capable of recognizing the audibility of objects and employ sounding object-aware separation, which helps the two models achieve significant better results. The experimental results can well validate the superiority of the proposed sounding object-aware separation mechanism.

We further show qualitative separation results for audible and silent objects in Fig. 5. We can see that both SoP and CoSep generate nonzero spectrograms for the silent *Guitar* and our CCoL can separate much better *Erhu* sound. The results can validate that the CCoL model is more robust to learning noise from silent objects during training and can effectively perform sounding object-aware sound separation.

Moreover, we train our model by letting the video 'B' and 'D' be randomly chosen from "with audio" and "silent".

Methods	SoP [49]	CoSep [13]	Random Obj	CoL	CCoL	GT
SDR	-11.35	-15.11	-11.34	14.78	91.07	264.44
SIR	-10.40	-12.81	-9.09	15.68	82.82	260.35

Table 3. Audio-visual sound separation performance (with silent objects). To help readers better appreciate the results, we include SDR and SIR from ground truth sounds.

In this way, each training video contains one or two sounding objects and the corresponding mixed sounds will be from up to four different sources. The SDR/SIR scores from SoP and CoL, and CCoL are 2.73/4.08, 5.79/11.43, and 6.46/11.72, respectively. The results further validate that the proposed cyclic co-learning framework can learn from both solo and duet videos and is superior over naive co-learning model.

From the sounding object visual grounding and audio-visual sound separation results, we can conclude that our cyclic co-learning framework can make the two tasks benefit from each other and significantly improve both visual grounding and sound separation performance.

4.4. Real-World Grounding and Separation

Besides synthetic data, our grounding and separation models can also handle real-world videos, in which multiple audible and silent objects might exist in single video frames. The example shown in Fig. 6 consists of three different instruments: *guitar*, *flute*, and *violin* in the video, in which *guitar* and *violin* are playing and making sounds. We can see that our grounding model can successfully identify the sounding objects: *guitar* and *violin* and silent object: *flute*. Meanwhile, our sounding object-aware separation model can separate individual sounds for each music instrument from the sound mixture. However, the SoP using the same noisy training data as our method obtains poor separation results because it is limited in associating objects

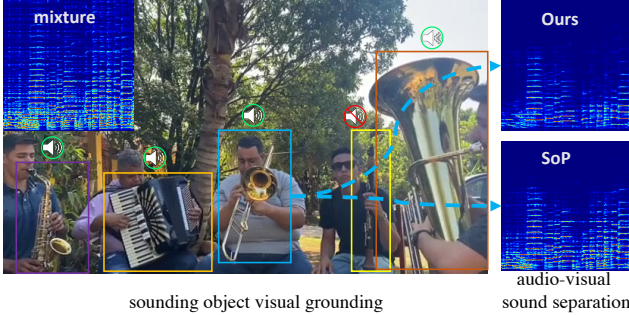


Figure 7. Real-world sounding object visual grounding and audio-visual sound separation for a challenging outdoor video with a 4-sound mixture and five instruments. From left to right, the instruments are *saxophone*, *accordion*, *trombone*, *clarinet* and *tuba*. Our grounding network can successfully find audible and silent objects, and our separation network can separate the *trombone* music, meanwhile suppressing unrelated sounds from other playing instruments. However, the separated sound by SoP contains much stronger noisy sounds from the *accordion* and *tuba*. Note that the *trombone* is an unseen category, and we have no 4-sound mixture during training. The example can demonstrate the generalization capacity of our method.

with corresponding sounds.

Another more challenging example is illustrated in Fig. 7. There are five different instruments in the video, and four of them are playing together. Although our grounding and separation networks are not trained with 4-sound mixtures, our sounding object visual grounding network can accurately find audible and silent objects in the video, and our audio-visual separation network can separate the *trombone* sound and is more capable of suppressing unrelated sounds from other three playing instruments than the SoP.

From the above results, we learn that our separation and grounding models trained using the synthetic data can be generalized to handle real-world challenging videos. In addition, the results further demonstrate that our models can effectively identify audible and silent objects and automatically discover the association between objects and individual sounds without relying on human annotations.

5. Limitation and Discussion

Our sounding object visual grounding model first processes video frames and obtains object candidates by a Faster R-CNN [36] object detector trained on a subset of Open Images dataset [30]. As stated in Sec. 3.2, such a strategy of using visual object proposal has been widely employed in image captioning [25, 2] and recent work on audio-visual learning [13]. Intuitively, similar to these previous works [25, 2, 13], the grounding model performance also highly relies on the quality of the object detector. If a sounding object is not detected by the detector, our model will fail to ground it. Thus, an accurate and effective object detector is very important for the grounding model. At

present, one possible way to address this limitation could be performing dense proposal of object candidates to prevent the above problem to some extent.

6. Conclusion and Future Work

In this paper, we introduce a cyclic co-learning framework that can jointly learn sounding object visual grounding and audio-visual sound separation. With the help of sounding object-aware sound separation to improve audio-visual sound separation. To further facilitate sounding object visual grounding learning, we use the separation model to help training sampling mining, which makes the learning process of the two tasks in a cycle and can simultaneously enhance both grounding and separation performance. Extensive experiments can validate that the two different problems are highly coherent, and they can benefit from each other with our cyclic co-learning, and the proposed model can achieve noticeable performance on both sounding object visual grounding and audio-visual sound separation.

There is various audio and visual content with different modality compositions in real-world videos. Besides the silent object issue, there are other challenging cases that affect audio-visual learning, as discussed below.

There might be multiple instances of the same instrument in a video. Our separation model mainly uses encoded visual semantic information to separate sounds; however, the visual dynamics that can provide additional discriminative information to identify different instances for the same instrument are not exploited. In the future, we will consider how to incorporate dynamic visual information to further strengthen our model’s ability on separating multiple sounds of the same types as in [48, 9].

Sound sources are not always visible in videos. For example, guitar sound might only be background music. In this case, there are no corresponding visual objects that can be used as conditions to separate the guitar sound for existing audio-visual sound separation methods. To address this problem, we might need to first parse input sound mixtures and video frames to recognize invisible sounds and then use other reliable conditions (*e.g.*, retrieved video frames or other semantic information) to separate the sounds.

Acknowledgement: Y. Tian and C. Xu were supported by NSF 1741472, 1813709, and 1909912. D. Hu was supported by Fundamental Research Funds for the Central Universities, the Research Funds of Renmin University of China (NO. 21XNLG17), the Beijing Outstanding Young Scientist Program NO. BJJWZYJH012019100020098 and the 2021 Tencent AI Lab Rhino-Bird Focused Research Program (NO. JR202141). The article solely reflects the opinions and conclusions of its authors but not the funding agents.

References

- [1] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In *ECCV*, 2020.
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.
- [3] Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *ECCV*, 2018.
- [4] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- [5] David A Bulkin and Jennifer M Groh. Seeing sounds: visual and auditory interactions in the brain. *Current opinion in neurobiology*, 16(4):415–419, 2006.
- [6] Lele Chen, Zhiheng Li, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Lip movements generation at a glance. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 520–535, 2018.
- [7] Lele Chen, Ross K Maddox, Zhiyao Duan, and Chenliang Xu. Hierarchical cross-modal talking face generation with dynamic pixel-wise loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7832–7841, 2019.
- [8] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation. *TOG*, 2018.
- [9] Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. Music gesture for visual sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10478–10487, 2020.
- [10] Chuang Gan, Hang Zhao, Peihao Chen, David Cox, and Antonio Torralba. Self-supervised moving vehicle tracking with stereo sound. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7053–7062, 2019.
- [11] Ruohan Gao, Rogerio Feris, and Kristen Grauman. Learning to separate object sounds by watching unlabeled video. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 35–53, 2018.
- [12] Ruohan Gao and Kristen Grauman. 2.5 d visual sound. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 324–333, 2019.
- [13] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019.
- [14] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. *arXiv preprint arXiv:1912.04487*, 2019.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [16] John R Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe. Deep clustering: Discriminative embeddings for segmentation and separation. In *ICASSP*, pages 31–35. IEEE, 2016.
- [17] John R Hershey and Javier R Movellan. Audio vision: Using audio-visual synchrony to locate sounds. In *Advances in neural information processing systems*, pages 813–819, 2000.
- [18] KAZUO Hikosaka, EIICHI Iwai, H Saito, and KEIJI Tanaka. Polysensory properties of neurons in the anterior bank of the caudal superior temporal sulcus of the macaque monkey. *Journal of neurophysiology*, 60(5):1615–1637, 1988.
- [19] Di Hu, Xuhong Li, Lichao Mou, P. Jin, D. Chen, L. Jing, X. Zhu, and D. Dou. Cross-task transfer for geotagged audiovisual aerial scene recognition. In *ECCV*, 2020.
- [20] Di Hu, Feiping Nie, and Xuelong Li. Deep multimodal clustering for unsupervised audiovisual learning. In *CVPR*, 2019.
- [21] Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou. Discriminative sounding objects localization via self-supervised audiovisual matching. In *Advances in Neural Information Processing Systems*, 2020.
- [22] Di Hu, Zheng Wang, Haoyi Xiong, Dong Wang, Feiping Nie, and Dejing Dou. Curriculum audiovisual learning. *arXiv preprint arXiv:2001.09414*, 2020.
- [23] Po-Sen Huang, Scott Deeann Chen, Paris Smaragdis, and Mark Hasegawa-Johnson. Singing-voice separation from monaural recordings using robust principal component analysis. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 57–60. IEEE, 2012.
- [24] Po-Sen Huang, Minje Kim, Mark Hasegawa-Johnson, and Paris Smaragdis. Joint optimization of masks and deep recurrent neural networks for monaural source separation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(12):2136–2147, 2015.
- [25] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [26] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5492–5501, 2019.
- [27] Einat Kidron, Yoav Y Schechner, and Michael Elad. Pixels that sound. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 88–95. IEEE, 2005.
- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [29] Bruno Korbar, Du Tran, and Lorenzo Torresani. Scsampler: Sampling salient clips from video for efficient action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6232–6242, 2019.

- [30] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*, 2017.
- [31] Yan-Bo Lin, Yu-Jhe Li, and Yu-Chiang Frank Wang. Dual-modality seq2seq network for audio-visual event localization. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2002–2006. IEEE, 2019.
- [32] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. *European Conference on Computer Vision (ECCV)*, 2018.
- [33] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, 2020.
- [34] Colin Raffel, Brian McFee, Eric J Humphrey, Justin Salamon, Oriol Nieto, Dawen Liang, Daniel PW Ellis, and C Colin Raffel. mir_eval: A transparent implementation of common mir metrics. In *ISMIR*, 2014.
- [35] Tanzila Rahman, Bicheng Xu, and Leonid Sigal. Watch, listen and tell: Multi-modal weakly supervised dense event captioning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8908–8917, 2019.
- [36] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.
- [37] Andrew Rouditchenko, Hang Zhao, Chuang Gan, Josh McDermott, and Antonio Torralba. Self-supervised segmentation and source separation on videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [38] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4358–4366, 2018.
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2014.
- [40] Barry E Stein and M Alex Meredith. *The merging of the senses*. The MIT Press, 1993.
- [41] Yapeng Tian, Chenxiao Guan, Goodman Justin, Marc Moore, and Chenliang Xu. Audio-visual interpretable and controllable video captioning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [42] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-supervised audio-visual video parsing. In *ECCV*, 2020.
- [43] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 247–263, 2018.
- [44] Xin Wang, Yuan-Fang Wang, and William Yang Wang. Watch, listen, and describe: Globally and locally aligned cross-modal attentions for video captioning. *arXiv preprint arXiv:1804.05448*, 2018.
- [45] Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention matching for audio-visual event localization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [46] Xudong Xu, Bo Dai, and Dahua Lin. Recursive visual sound separation using minus-plus net. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 882–891, 2019.
- [47] Xudong Xu, Hang Zhou, Ziwei Liu, Bo Dai, Xiaogang Wang, and Dahua Lin. Visually informed binaural audio generation without binaural audios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [48] Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1735–1744, 2019.
- [49] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018.
- [50] Hang Zhou, Yasheng Sun, Wu Wayne, Chen Change Loy, Xiaogang Wang, and Liu Ziwei. Pose-controllable talking face generation by implicitly modularized audio-visual representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [51] Hang Zhou, Xudong Xu, Dahua Lin, Xiaogang Wang, and Ziwei Liu. Sep-stereo: Visually guided stereophonic audio generation by associating source separation. In *European Conference on Computer Vision*, pages 52–69. Springer, 2020.