

Positive Sample Propagation along the Audio-Visual Event Line

Jinxing Zhou^{1,2} Liang Zheng^{3*} Yiran Zhong³ Shijie Hao^{1,2} Meng Wang^{1,2*}

¹Hefei University of Technology

²Intelligent Interconnected Systems Laboratory of Anhui Province

³Australian National University

{zhoujxhfut, hfut.hsj, eric.mengwang}@gmail.com,

{liang.zheng, yiran.zhong}@anu.edu.au

Abstract

Visual and audio signals often coexist in natural environments, forming audio-visual events (AVEs). Given a video, we aim to localize video segments containing an AVE and identify its category. In order to learn discriminative features for a classifier, it is pivotal to identify the helpful (or positive) audio-visual segment pairs while filtering out the irrelevant ones, regardless whether they are synchronized or not. To this end, we propose a new positive sample propagation (PSP) module to discover and exploit the closely related audio-visual pairs by evaluating the relationship within every possible pair. It can be done by constructing an all-pair similarity map between each audio and visual segment, and only aggregating the features from the pairs with high similarity scores. To encourage the network to extract high correlated features for positive samples, a new audio-visual pair similarity loss is proposed. We also propose a new weighting branch to better exploit the temporal correlations in weakly supervised setting. We perform extensive experiments on the public AVE dataset and achieve new state-of-the-art accuracy in both fully and weakly supervised settings, thus verifying the effectiveness of our method.

1. Introduction

Recent literature has shown that by fusing multi-modality information can lead to better deep feature presentation, *i.e.*, audio-visual fusion [2] and text-visual fusion [20]. However, building a large scale multi-modality pre-training datasets would require heavy manual labours to clean and annotate the raw video sets. To relief the manual labour, recent work either focuses on learning from noise supervision [5, 12] or tries to automatically filter out unpaired samples [28].

The task of Audio-Visual Event (AVE) localization [28]



Figure 1. An illustration of the AVE localization task. Each video segment is composed of an audio and a visual component. In this example, the “hum” of the bus exists in all the segments (audio modality), but the visual images of the “bus” only appear in the third and fourth segments (visual modality). So only these two segments (red boxes) are localized as (bus) event, the remaining are recognized as background.

is served for the latter purpose. An AVE often refers as an event that is both audible and visible in a video segment, *i.e.*, a sound source appears in an image (*visible*) while the source of the sound also exists in audio portion (*audible*). As shown in Fig. 1, a bus humming is an AVE in the third and fourth segments as we can see a bus and hear it humming simultaneously in these video segments. The AVE localization task is to find these video segments that contain an audio-visual event and classify it into a certain category¹.

There are two relations that need to be considered in the AVE task: intra-modal relations and cross-modal relations. The former often addresses temporal relations in one single modality while the later also takes audio and visual relations into account. The pioneer work [14, 28] often tries to regress the class by concatenating features from synchronized audio-visual pairs. Since these methods do not explicitly consider the intra-modal or cross-modal relations, their accuracy is often unsatisfying. The follow-

¹Note that there is a fundamental difference between the Multimedia Event Detection (MED) task and the AVE localization task: MED is a retrieval task that aims to find video clips that are associated with a particular event from a video archive while AVE localization is a classification problem.

*Corresponding author.

ing works [27, 30, 31, 32] utilize a self-attention mechanism to explicitly encode the temporal relations within intra-modality and some of them [21, 22, 31, 32] also aggregate better audio-visual feature representations by encoding cross-modal relations. However, these methods often ignore the interference caused by irrelevant audio-visual segment pairs during the fusion process. In this paper, we argue that by only aggregating features from positive samples, *i.e.*, high-relevant audio-visual pairs, we can have better AVE localization accuracy.

Specifically, we propose a new Positive Sample Propagation (PSP) module. In a nutshell, PSP first constructs an all-pair similarity map between each audio and visual segment and cuts off the entries that are below a pre-set similarity threshold, and then aggregates the audio and visual features without considering the negative and weak entries in an online fashion. Through various visualizations we show that PSP allows more relevant features that are not necessarily synchronized to be aggregated in an online fashion.

Apart from PSP that can be used in both fully and weakly supervised settings (refer Sec. 3 for the setting details), we further propose two improvements that work under each setting, respectively. On the one hand, an audio-visual pair similarity loss is introduced under the fully supervised setting that encourages the network to learn high correlated features of audio and visual segments if they belong to the same event. On the other hand, we propose a weighting branch in the weakly supervised setting, which gives temporal weights to the segment features.

We evaluate our method on the standard AVE dataset [28]. We show that the proposed techniques consistently benefit our system and when combined allow us to achieve state-of-the-art performance under both fully and weakly supervised settings.

2. Related work

Audio-visual correspondence (AVC) aims to predict whether a given visual image corresponds duration of the audio. A model is asked to judge whether the audio and visual signals describe the same object, *e.g.*, dog *v.s.* bark, cat *v.s.* meow. It is a self-supervised problem since the visual image is usually accompanied by the corresponding sound. Existing methods try to evaluate the correspondences by measuring the audio-visual similarity [2, 3, 4, 6, 8, 11]. It will get a large similarity score if the audio-visual pair is corresponding, otherwise, a low score. This motivates us to tackle the abundant audio-visual pairs in the AVE localization problem by considering the audio-visual similarity.

Sound source localization aims to localize those visual regions which are relevant to the provided audio signal. It is related to *sound source separation* [1, 7, 17, 18, 33] problem. The target region of the visual frame must be corresponding with the given sound. It is similar to the AVC

task from this point of view, but the real challenge of sound source localization is to accurately locate the sound-maker when there are multiple sound sources in a visual frame. Qian *et al.* [19] adapt the Grad-CAM [24] to disentangle class-specific features for multiple sound sources problem. Senocak *et al.* [25] propose a triplet loss working in an unsupervised manner. Afouras *et al.* [1] utilize a contrastive loss to train the model in a self-supervised learning way. Both of these methods [1, 25] need to construct positive and negative audio-visual pair samples. Since similar positive and negative samples are easily obtained in AVE localization, depending on whether the audio and visual segments depict the same event, we try to research those audio-visual pairs and explore its effect.

Audio-visual event localization aims to distinguish those segments including an audio-visual event from a long video. Existing works mainly focus on the audio-visual fusion process. A dual multimodal residual network is proposed in [28]. Lin *et al.* [14] adapt a bi-directional LSTM [23] to fuse audio and visual features in a seq2seq manner. During the whole fusion process, simple concatenation and addition operations are adapted along the single synchronized audio-visual pair. Ramaswamy [21, 22] utilizes a bilinear method to capture cross-modal relations. Xuan *et al.* [32] propose to leverage *modality sentinel* to give different weights to audio and visual features. Lin *et al.* [15] design an audio-visual transformer to describe local spatial and temporal information. The visual frame is divided into patches and adjacent frames are utilized, making the model complicated and computationally intensive. Xu *et al.* [31] attempt to leverage concatenating audio-visual features as the supervision then the feature of each modality is updated by separate modules. Unlike these, the proposed PSP method has a further in-depth study on the abundant audio-visual pairs, selecting the most relevant ones. Relying on these positive samples, more distinguished audio-visual features can be obtained after feature aggregation.

3. Problem statement

AVE localization aims to find out those segments containing an audio-visual event [28]. In other words, AVE localization is expected to decide whether each synchronized audio-visual pair depicts an event. Besides, AVE localization needs to identify the event category for each segment. Specifically, a video sequence S is divided into T non-overlapping yet continuous segments $\{S_t^v, S_t^a\}_{t=1}^T$, and each segment is one-second long. S^v and S^a are the visual and audio components, respectively. We consider two settings of this task, to be described below.

Fully-supervised AVE localization. Under the fully-supervised setting, the event label of every video segment is given, indicating whether the segment denotes an event and which category the event belongs to. We

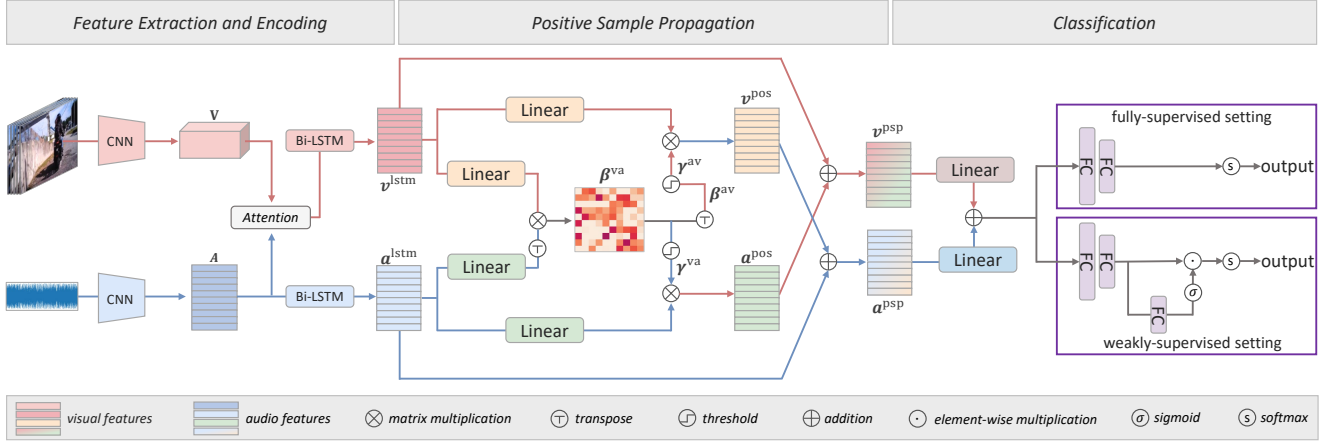


Figure 2. System Flow. We first extract and encode video and audio features through existing modules such as AVGA [28] and Bi-LSTM. The proposed positive sample propagation (PSP) takes the LSTM encoded features as input, which are fed to a few linear layers. An affinity matrix is computed before selecting the positive connections of audio-visual segment pairs using thresholding. In this module, audio and visual features are aggregated by feature propagation through the positive connections. In the last stage, we classify the event into predefined categories. For the supervised setting, apart from the commonly used CE loss, we further propose an audio-visual pair similarity loss which enforces similar features between them when they contain an event. For the weakly supervised setting, we introduce another FC layer that gives weights to different video segments: higher weights are given event-containing segments.

denote the event label of the t^{th} segment as $\mathbf{y}_t = \{y_t^c | y_t^c \in \{0, 1\}, \sum_{c=1}^C y_t^c = 1\} \in \mathbb{R}^C$, where C is the number of categories (including the *background*). Then, the label for the entire video can be written as $\mathbf{Y}^{\text{fully}} = [\mathbf{y}_1; \mathbf{y}_2; \dots; \mathbf{y}_T] \in \mathbb{R}^{T \times C}$. Through $\mathbf{Y}^{\text{fully}}$, we know whether an arbitrary synchronized audio-visual pair at time t is an event: if the 1 of its event label \mathbf{y}_t is at the entry of a certain event instead of the *background*, the pair describes an event, and otherwise does not.

Weakly-supervised AVE localization. We adapt the weakly-supervised setting proposed in [14, 32], where the label $\mathbf{Y}^{\text{weak}} \in \mathbb{R}^{1 \times C}$ is the average pooling value of $\mathbf{Y}^{\text{fully}}$ along the column. It implies the proportion of audio-visual pairs that contain an event. This setting is different from the fully-supervised one because the event label of each segment y_t is unknown, making the problem more challenging.

4. Our method

4.1. Overall pipeline

The overall pipeline of our system is illustrated in Fig. 2, which includes three modules: a feature extraction and encoding module, a positive sample propagation module, and a classification module. In the *feature extraction and encoding* module, audio-guided visual attention (AVGA [28]) is adapted for early fusion to make the model focus on those visual regions closely related to the audio component. Then a Bi-LSTM is utilized to encode temporal relations in video segments. The LSTM encoded features are sent to the proposed *positive sample propagation (PSP)* module. PSP is able to select those positive connections of audio-visual segment pairs by measuring the cross-modal similarity with thresholding. Audio and visual features are aggregated by

feature propagation through the positive connections. The updated audio-visual features after PSP are fused then sent to the final *classification* module, predicting which video segments contain an event and the event category.

4.2. Feature extraction and encoding

The visual and synchronized audio segments are processed by pretrained convolutional neural networks (CNNs). We denote the resulting visual feature as $\mathbf{V} \in \mathbb{R}^{T \times N \times d_v}$, where d_v is the feature dimension, $N = H \times W$, H and W are the height and width of the feature map, respectively. The extracted audio feature is denoted as $\mathbf{A} \in \mathbb{R}^{T \times d_a}$, where d_a denotes feature dimension. We then directly adapt AVGA [28] for multi-modal early fusion. AVGA allows the model to focus on visual regions that are relevant to the audio component. To encode the temporal relationship in video sequences, the visual and audio features after AVGA are further sent to two independent Bi-LSTMs. The updated visual and audio features are represented as $\mathbf{v}^{\text{lstm}} \in \mathbb{R}^{T \times d_l}$ and $\mathbf{a}^{\text{lstm}} \in \mathbb{R}^{T \times d_l}$, respectively.

4.3. Positive sample propagation (PSP)

PSP allows the network to learn more representative features by exploiting the similarities of audio-visual pairs. It involves three steps.

In *all-pair connection construction*, all the audio-visual pairs are connected. As shown in Fig. 3, here we only display the connections of one visual segment for simplicity, i.e., $\langle v_1 \leftrightarrow a_1/a_2/a_3/a_4 \rangle$. The strength of these connections are measured by the similarity between the audio-visual components $\langle \mathbf{a}^{\text{lstm}}, \mathbf{v}^{\text{lstm}} \rangle$, computed by,

$$\beta^{\text{va}} = \frac{(\mathbf{v}^{\text{lstm}} \mathbf{W}_1^v)(\mathbf{a}^{\text{lstm}} \mathbf{W}_1^a)^\top}{\sqrt{d_l}}, \quad \beta^{\text{av}} = (\beta^{\text{va}})^\top, \quad (1)$$

where \mathbf{W}_1^v and $\mathbf{W}_1^a \in \mathbb{R}^{d_l \times d_h}$ are learnable parameters of linear transformations, implemented by a linear layer, and d_l is the dimension of the audio or visual feature. β^{va} and $\beta^{av} \in \mathbb{R}^{T \times T}$ are the similarity matrices.

Second, we *prune the negative and weak connections*. Specifically, the connections constructed in the first step are divided into three groups according to the similarity values: negative, weak, and positive. As a classification task, the success of AVE localization highly depends on the richness and correctness of training samples for each class. That is, we aim to collect possibly many and relevant *positive* connections. We achieve this goal by filtering out the weak and negative ones, *e.g.*, $v_1 \leftrightarrow a_3$ and $v_1 \leftrightarrow a_4$ as shown in Fig. 3. We begin with processing all the audio-visual pairs with the ReLU activation function, cutting off connections with negative similarity values. Row-wise ℓ_1 normalization is then performed, yielding the normalized similarity matrices β^{va} and β^{av} .

The negative and weak connections are presumably featured by smaller similarity values, so we simply adapt a thresholding method, written as,

$$\begin{aligned} \gamma^{va} &= \beta^{va} \mathbb{I}(\beta^{va} - \tau), \\ \gamma^{av} &= \beta^{av} \mathbb{I}(\beta^{av} - \tau), \end{aligned} \quad (2)$$

where τ is the hyper-parameter, controlling how many connections will be pruned. $\mathbb{I}(\cdot)$ is an indicator function, which outputs 1 when the input is greater than or equal to 0, and otherwise outputs 0, γ^{va} . After thresholding, row-wise ℓ_1 normalization is again performed to obtain the final similarity matrices $\gamma^{va}, \gamma^{av} \in \mathbb{R}^{T \times T}$.

Online feature aggregation. The above step identifies audio (visual) components with high similarities with a given visual (audio) component, *e.g.*, $v_1 \leftrightarrow a_1$ and $v_1 \leftrightarrow a_2$ shown in Fig. 3. This is essentially a positive sample propagation process that can be utilized to update the features of audio or visual components. Particularly, given the connection weights γ^{av} and γ^{va} , the audio and visual features \mathbf{a}^{PSP} and \mathbf{v}^{PSP} are respectively updated as,

$$\begin{aligned} \mathbf{a}^{\text{PSP}} &= \overbrace{\gamma^{av} (\mathbf{v}^{\text{Istm}} \mathbf{W}_2^v)}^{\mathbf{v}^{\text{pos}}} + \mathbf{a}^{\text{Istm}}, \\ \mathbf{v}^{\text{PSP}} &= \overbrace{\gamma^{va} (\mathbf{a}^{\text{Istm}} \mathbf{W}_2^a)}^{\mathbf{a}^{\text{pos}}} + \mathbf{v}^{\text{Istm}}, \end{aligned} \quad (3)$$

where $\mathbf{W}_2^a, \mathbf{W}_2^v \in \mathbb{R}^{d_l \times d_l}$ are parameters defining linear transformations, and $\mathbf{a}^{\text{PSP}}, \mathbf{v}^{\text{PSP}} \in \mathbb{R}^{T \times d_l}$.

Generally, the audio (visual) feature \mathbf{a}^{PSP} (\mathbf{v}^{PSP}) is enhanced by the propagated positive support from the other modality. This practice allows us to learn more discriminative audio-visual representations, displayed in Fig. 5. More discussions are provided in Sec. 4.6.

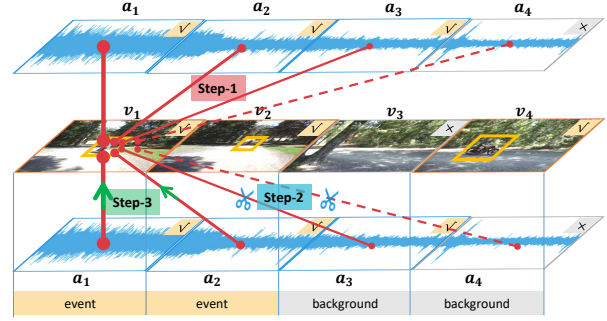


Figure 3. An illustration of the proposed PSP. In this example, only the first two video segments contain an audio-visual event, *i.e.*, motorcycle. “ \surd ” denotes the audio or visual segment describes the event, while “ \times ” means not. The red lines denote connections of audio-visual pairs, solid lines represent connections formed by relevant pairs, while dotted lines denote irrelevant pairs. The thickness of line reflects the similarity of the audio-visual pair. $v_1 \leftrightarrow a_4$ is a *negative* connection, formed by irrelevant audio-visual pair with negative similarity value. $v_1 \leftrightarrow a_3$ and $v_1 \leftrightarrow a_1/a_2$ are *weak* and *positive connections* respectively, determined via similarity. The upper part corresponds to “Step-1” (all-pair connection construction), while the lower part denotes “Step-2” (prune the negative and weak connections), and the green arrow indicates the direction of feature propagation (“Step-3”).

4.4. Classification

Before classifier prediction, we transform the visual and audio features into the same embedding space through another linear layer, and then combine the output through simple averaging, yielding the fusion feature denoted as $\mathbf{f}^{v \leftrightarrow a}$. This process is written as,

$$\mathbf{f}^{v \leftrightarrow a} = \frac{1}{2} [\mathcal{N}(\mathbf{v}^{\text{PSP}} \mathbf{W}_3^v) + \mathcal{N}(\mathbf{a}^{\text{PSP}} \mathbf{W}_3^a)], \quad (4)$$

where $\mathcal{N}(\cdot)$ represents layer normalization, $\mathbf{W}_3^v, \mathbf{W}_3^a \in \mathbb{R}^{d_l \times d_l}$ represent learnable parameters in the linear layers, and $\mathbf{f}^{v \leftrightarrow a} \in \mathbb{R}^{T \times d_l}$.

For the fully-supervised setting, as shown in Fig. 2, the fusion feature is further processed by two FC layers. The classifier prediction $\mathbf{o}^{\text{fully}} \in \mathbb{R}^{T \times C}$ can be obtained through a softmax function.

For the weakly supervised setting, different from existing methods [14, 28, 32], we add a weighting branch on the fully supervised classification module (Fig. 2). It is essentially another FC layer that enables the model to further capture the differences between synchronized audio-visual pairs by dynamically focusing on different event categories. This process is summarized below,

$$\begin{aligned} \mathbf{f}^h &= \mathbf{f}^{v \leftrightarrow a} \mathbf{W}_4^{\text{weak}} \mathbf{W}_5^{\text{weak}}, \\ \phi &= \sigma(\mathbf{f}^h \mathbf{W}_6^{\text{weak}}), \\ \mathbf{o}^{\text{weak}} &= s(f_{\text{avg}}(\mathbf{f}^h \odot \Phi)), \end{aligned} \quad (5)$$

where $\mathbf{W}_4^{\text{weak}} \in \mathbb{R}^{d_l \times d_h}$, $\mathbf{W}_5^{\text{weak}} \in \mathbb{R}^{d_h \times C}$, $\mathbf{W}_6^{\text{weak}} \in \mathbb{R}^{C \times 1}$ are learnable parameters in the FC layers, and $\mathbf{f}^h \in$

$\mathbb{R}^{T \times C}$. σ and s denote the sigmoid and softmax operators, respectively. $\phi \in \mathbb{R}^{T \times 1}$ weighs the importance of the temporal video segments, and $\Phi \in \mathbb{R}^{T \times C}$ is obtained by duplicating ϕ for C times. \odot is the element-wise multiplication, f_{avg} is the average operation along the temporal dimension. The final prediction $\mathbf{o}^{\text{weak}} \in \mathbb{R}^{1 \times C}$.

For comparison, we denote predictions through a network without the weighting branch as,

$$\mathbf{o}_{\text{wo}}^{\text{weak}} = s(f_{\text{avg}}(\mathbf{f}^h)). \quad (6)$$

4.5. Objective function

Fully supervised setting. Given the network output $\mathbf{o}^{\text{fully}}$ and ground truth $\mathbf{Y}^{\text{fully}}$, we adapt the cross entropy (CE) loss as the objective function, written as,

$$\mathcal{L}_{\text{ce}} = -\frac{1}{TC} \sum_{t=1}^T \sum_{c=1}^C \mathbf{Y}_{tc}^{\text{fully}} \log(\mathbf{O}_{tc}^{\text{fully}}) \quad (7)$$

Recall that each row of $\mathbf{Y}^{\text{fully}}$ contains a one-hot event label vector, describing the category of the corresponding segment (synchronized audio-visual pair). As such, this classification loss allows the network to predict which *event category* a video segment contains.

Apart from the CE loss, we propose a new loss item, named audio-visual pair similarity loss $\mathcal{L}_{\text{avps}}$. In principle, it asks the network to produce similar features for a pair of audio and visual components if the pair *contains an event*. Specifically, for a video composed of T segments, we define label vector $\mathbf{G} = \{g_t | g_t \in \{0, 1\}, t = 1, 2, \dots, T\} \in \mathbb{R}^{1 \times T}$, where g_t represents whether the t^{th} segment is an event or background. Next, ℓ_1 normalization is performed on \mathbf{G} . We then compute the ℓ_1 normalized similarity vector $\mathbf{S} \in \mathbb{R}^{1 \times T}$ between the visual and audio features

$$\mathbf{S} = \frac{\mathbf{v}^{\text{psp}} \odot \mathbf{a}^{\text{psp}}}{\|\mathbf{v}^{\text{psp}} \odot \mathbf{a}^{\text{psp}}\|_1}, \quad (8)$$

where $\|\cdot\|_1$ calculate the ℓ_1 norm of a vector. The proposed loss $\mathcal{L}_{\text{avps}}$ is then written as,

$$\mathcal{L}_{\text{avps}} = \mathcal{L}_{\text{MSE}}(\mathbf{S}, \mathbf{G}), \quad (9)$$

where $\mathcal{L}_{\text{MSE}}(\cdot, \cdot)$ computes the mean squares error between two vectors.

Combining Eq. 9 and Eq. 7, the overall objective function for fully-supervised setting $\mathcal{L}_{\text{fully}}$ can be computed by:

$$\mathcal{L}_{\text{fully}} = \mathcal{L}_{\text{ce}} + \lambda \mathcal{L}_{\text{avps}}, \quad (10)$$

where λ is a hyper-parameter to balance the two losses.

Weakly supervised setting. For this setting, following the practice in [14, 31], we adapt the binary cross entropy (BCE) loss, formulated as,

$$\begin{aligned} \mathcal{L}_{\text{w-bce}} &= \mathcal{L}_{\text{BCE}}(\mathbf{o}^{\text{weak}}, \mathbf{Y}^{\text{weak}}), \\ \mathcal{L}_{\text{wo-bce}} &= \mathcal{L}_{\text{BCE}}(\mathbf{o}_{\text{wo}}^{\text{weak}}, \mathbf{Y}^{\text{weak}}), \end{aligned} \quad (11)$$

where $\mathcal{L}_{\text{w-bce}}$ and $\mathcal{L}_{\text{wo-bce}}$ are calculated between the ground-truths and predictions. \mathbf{o}^{weak} and $\mathbf{o}_{\text{wo}}^{\text{weak}}$ are predictions obtained with or without the weighting branch (Sec. 4.4 and Fig. 2), respectively.

4.6. Discussion

Detailed examination and meanings of \mathbf{v}^{pos} and \mathbf{a}^{pos} .

The computation of \mathbf{v}^{pos} (\mathbf{a}^{pos}) is shown in Eq. 3. Take \mathbf{v}^{pos} for example. The i^{th} row $\mathbf{v}_i^{\text{pos}}$ is the weighted sum of the visual feature $\mathbf{v}_j^{\text{lstm}}$ ($j = 1, 2, \dots, T$) after linear transformation. Here the weight, denoted as γ_i^{av} , is exactly the similarity between the audio feature \mathbf{a}_i and features of all the visual components. Note that some elements of γ_i^{av} are zeros since the negative and weak connections are pruned during PSP, so $\mathbf{v}_i^{\text{pos}}$ is the aggregation result of those *positive* visual features which are most relevant to \mathbf{a}_i .

Physical meanings of \mathbf{v}^{psp} and \mathbf{a}^{psp} . Take \mathbf{a}^{psp} for example. From Eq. 3, we find \mathbf{a}^{psp} is composed of two features: the original audio feature \mathbf{a}^{lstm} and the aggregation of positive visual features \mathbf{v}^{pos} . As discussed above, those positive visual features have large audio-visual similarity values, *i.e.*, small vector angles and similar vector directions. Therefore, after being added to \mathbf{v}^{pos} , the magnitude and direction of vectors representing original audio feature \mathbf{a}^{lstm} will be changed to reflect that during training. Such an adjustment in the distribution of audio representation can be verified by the visualization results in Fig. 5.

Why an additional FC layer in the weakly supervised setting? When fully supervised, clear supervision is known for each segment. For the weakly supervised setting, both the ground truth label $\mathbf{Y}^{\text{weak}} \in \mathbb{R}^{1 \times C}$ and the prediction $\mathbf{o}^{\text{weak}} \in \mathbb{R}^{1 \times C}$ are obtained through an average pooling operation along the temporal dimension. Without knowing the supervision for each segment, the baseline approach considers all temporal video segments to have similar weights when calculating the loss. It makes it harder for the model to focus on video segments that contain an event. In our design, through the sigmoid activation function, we obtain the weights of temporal video segments. As such, our model can better distinguish these temporal sequences and thus help locate which segments contain an event.

Implications of Eq. 9. As shown in Eq. 7, the classification loss \mathcal{L}_{ce} prompts the model to correctly predict the event *categories*. In comparison, $\mathcal{L}_{\text{avps}}$ allows the network to be aware of *whether an event exists in an audio-visual pair*. Specifically, if g_t is equal to 1, the synchronized audio-visual feature should have a higher similarity, and otherwise lower. Therefore, for an audio (visual) component, $\mathcal{L}_{\text{avps}}$ provides another auxiliary constraint so that the model can better select the most relevant visual (audio) components for feature aggregation during PSP. Note that $\mathcal{L}_{\text{avps}}$ cannot be adapted in the weakly supervised setting, where the label g_t of each segment is unknown.

5. Experiment

5.1. Experimental setup

Dataset. Following existing works [14, 28, 31, 32], we use the AVE dataset [28] that is publicly available. This dataset contains 4,143 videos, which cover various real-life scenes and can be divided into 28 event categories, *e.g.*, church bell, male speech, acoustic guitar, and dog barking. Each video sample is evenly partitioned into 10 segments, and each segment is one-second long. The audio-visual event boundary on the segment level and the event category on the video level are provided.

Evaluation metric. The category label of each segment is predicted in both fully and weakly supervised settings. Following [14, 28, 31, 32], we adapt the classification accuracy of each segment as the evaluation metric.

Implementation details. We use VGG-19 [26] pretrained on ImageNet [13] to extract the visual features. Specifically, 16 frames are sampled from each one-second segment. We extract the visual feature maps from each frame and use the average map as the visual feature for this segment. For audio features, we first process the raw audio into log-mel spectrograms and then extract the acoustic features using a VGG-like network [10] pretrained on AudioSet [9]. Besides, dropout technique is used in the linear layers (Fig. 2). Weight λ in Eq. 10 is empirically set to 100.

5.2. Quantitative Analysis

The effectiveness of the PSP encoding can be verified through comparing with an ablation study, *i.e.*, removing it from the localization network (Fig. 2). In Table 1, We denote the method without PSP as “w/o PSP”. We observe from the table that the performance drops in both the fully supervised and weakly supervised setting significantly. Specifically, the accuracy decrease is 4.1% (from 77.8% to 73.7%) and 3.3% (from 73.5% to 70.2%) for the two settings, respectively. This experiment clearly validates PSP.

Comparison with alternative positive sample selection methods. In our method, we emphasize that weak and negative samples are filtered out. Here, we compare this strategy with two variants: 1) all connections are used; 2) only negative ones are removed. Results are shown in Table 1. We have two main observations.

First, when all samples are propagated (denoted as “ASP”), the accuracy drops by 1.9% and 2.3% on the two settings, respectively. This shows that it is essential to have a selection process before feature aggregation instead of utilizing all the connections. In fact, the ASP variant shares the same spirit with HAN [27].

Second, when we only remove the negative connections (*i.e.*, those with a similarity value below $\tau = 0$), the system is inferior to the full method. Specifically, the classification accuracy decreases by 1.8% and 2.3% under the two set-

Method	Fully-supervised	Weakly-supervised
w/o PSP	73.7	70.2
ASP	75.9	71.2
WPSP	76.0	71.2
SAPSP	75.4	70.8
PSP (ours)	77.8	73.5

Table 1. Ablation studies of the proposed PSP, measured by accuracy(%) on the AVE dataset. “w/o” denotes “without”. “ASP” means retaining all connections ($\tau = -\infty$), while “WPSP” uses the weak and positive ones ($\tau = 0$). “SAPSP” represents adding self-attention to the feature extractor.

tings, which validates the effectiveness of filtering out the negative connections.

Comparison with adding self-attention [29] to the feature extractor. Self-attention [29] is widely used in existing methods [27, 30, 31, 32] to capture relationships within single modality. To explore whether it is useful in our system, we add a self-attention module before the Bi-LSTMs and denote it as the “SAPSP” method. As shown in Table 1, the performance surprisingly decreases by 2.4% and 2.7% under fully and weakly supervised settings, respectively. We speculate that the PSP module is sufficient to describe the cross-modality while implicitly reveals the intra-modality correlations. For example, in PSP a visual component is constrained to have similar features with multiple audio components describing the same event. Such cross-modality similarity at the same time implies that the similarity of the involved audio components to be high. In our future work, we will study in-depth the intra-modality and inter-modality similarities.

Benefit of the audio-visual pair similarity loss \mathcal{L}_{avps} . We respectively adapt \mathcal{L}_{ce} and $\mathcal{L}_{ce} + \lambda\mathcal{L}_{avps}$ as the objective function for model training. Two baselines are used: our PSP system and the AVEL system [28]. Results are presented in Table 2. We can clearly see that \mathcal{L}_{avps} improves the accuracy when the system is fully supervised. The improvement is 1.2% and 1.5% for PSP and AVEL, respectively. These results verify the role of \mathcal{L}_{avps} as an auxiliary restriction to help to select the positive audio-visual pairs for feature aggregation.

Improvement from the additional FC in the weakly supervised setting. In the weakly supervised setting, the major difference between our classification module and traditional methods [14, 28, 32] consists in the weighting branch (Fig. 2). To evaluate its effectiveness, we also implement this branch on top of the PSP and AVEL baselines. The results are shown in the last two rows of Table 2. We find that the performance of PSP and AVEL is improved by 1.9% and 2.3%, respectively. We argue the additional weighting branch within the designed classification module allows the model to give different weights to the temporal sequences, thus benefiting the localization of the target video segments. These results confirm the effectiveness of

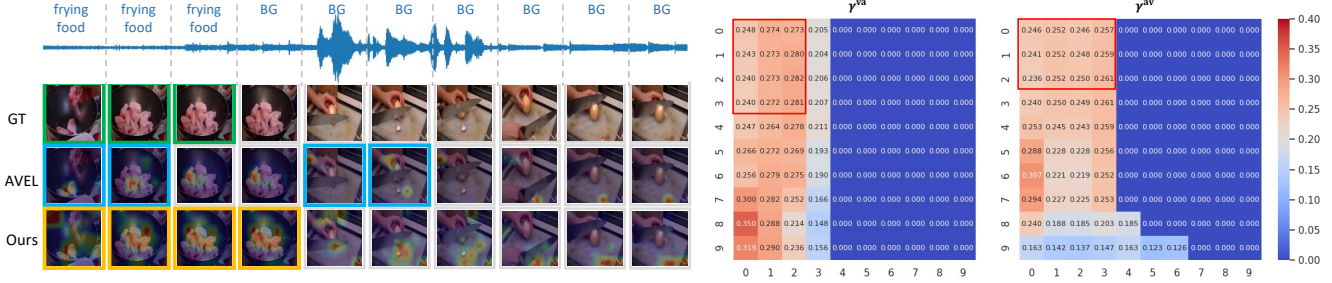


Figure 4. A qualitative example of AVE localization. For the video on the left, only the first three segments contain the visual and audio signals of the event *frying food*. The green boxes represent ground truth labels. The blue and orange boxes indicate predictions of AVEL [28] and our method, respectively. Besides, we visualize the attention effect on the images. It is clear that our method produces more accurate localization and that our attended regions better overlap with the sound sources. On the right, we visualize the audio-visual similarity matrices γ^{va} and γ^{av} (Eq. 3) after PSP. For γ^{va} , the x-axis and y-axis correspond to audio and visual features, respectively, and for γ^{av} the order is reversed. The red bounding boxes in γ^{va} show that the first three audio components are highly correlated with the first four visual components. Besides, negative and weak connections are cut off to 0 in γ^{va} and γ^{av} . The color bar corresponds to the similarity strength, with red denoting high similarities and blue for low similarities.

Setting	Method	PSP (ours)	AVEL [28]
fully	\mathcal{L}_{ce}	76.6	69.8*
	$\mathcal{L}_{ce} + \lambda\mathcal{L}_{avps}$	77.8	71.3*
weakly	w/o weight. branch	71.6	66.9*
	w/ weight. branch	73.5	69.2*

Table 2. Method comparison on the AVE dataset under two settings. We evaluate 1) the audio-visual pair similarity loss \mathcal{L}_{avps} under the fully supervised setting, and 2) the weighting branch under the weakly supervised setting. The two improvements are implemented on top of our system and AVEL [28]. Under AVEL, * denotes that the number is produced by us. We use **bold** font to show the higher performance brought by our technique.

τ	0	0.025	0.075	0.095	0.115
Fully-supervised	75.9	76.1	75.3	77.8	76.6
Weakly-supervised	71.2	71.7	70.4	73.5	72.8

Table 3. Impact of various values of τ on the system accuracy. Results on the two setting are shown.

the proposed improvements. We refer readers to Sec. 4.6 for discussions on the two techniques.

Sensitivity to hyper-parameter τ . The selection process is controlled by τ , determining how many connections will be cut off. Its influence on the system accuracy is shown in Table 3. We observe that overall the accuracy remains stable when τ varies between 0 and 0.115 and that the highest accuracy is achieved when $\tau = 0.095$. For different videos, the proportion of segments that are cut off highly depends on the video itself. If the whole video contains the same event of interest, it is likely that most will be retained in training; if a video contains lots of background, the same threshold will cut off more of its content.

Comparison with the state of the art. We compare our method with the state of the art in Table 4, where we report superior results: **the classification accuracy is 77.8% and 73.5% for the fully and weakly supervised settings, respectively.** Compared with the baseline feature extrac-

Method	Fully-supervised	Weakly-supervised
AVEL [28]	68.6	66.7
AVSDN [14]	72.6	67.3
CMAN [32]	73.3*	70.4*
DAM [30]	74.5	-
AVRB [22]	74.8	68.9
AVIN [21]	75.2	69.4
AVT [15]	76.8	70.2
CMRA [31]	77.4	72.9
PSP (Ours)	77.8	73.5

Table 4. Comparison with the state-of-the-art methods under two settings, measured by accuracy(%) on the AVE dataset. * indicates the number is reproduced by us.

tor AVEL [28], we exceed it by 9.2% and 6.8% under the fully and weakly supervised settings, respectively. This can also be proved by the results shown in Table 2 where the numbers of AVEL are reproduced by ourselves. Moreover, while the AVGA module [28] adapted in our system is slightly lower (0.6%) than the recent audio-guided spatial-channel attention (AGSCA) [31], our overall system manages to obtain higher accuracy than AGSCA. This can be attributed to both the PSP and our system design.

5.3. Qualitative analysis

We start by presenting an example of audio-visual event localization in Fig. 4. The event in this sample is difficult to predict because the visual images are changeable and the audio signals are mixed with background noise. 1). While both our method and AVEL [28] use the AVGA attention, we show that our method enables better attention to visual regions closely related to sound sources. As displayed in Fig. 4, for the event of *frying food*, our attended regions include both the frying chicken thighs and the pot, especially in the first four segments. In comparison, AVEL only finds the thighs and very small receptive fields. 2). Our method has a better prediction result. AVEL seems to make deci-

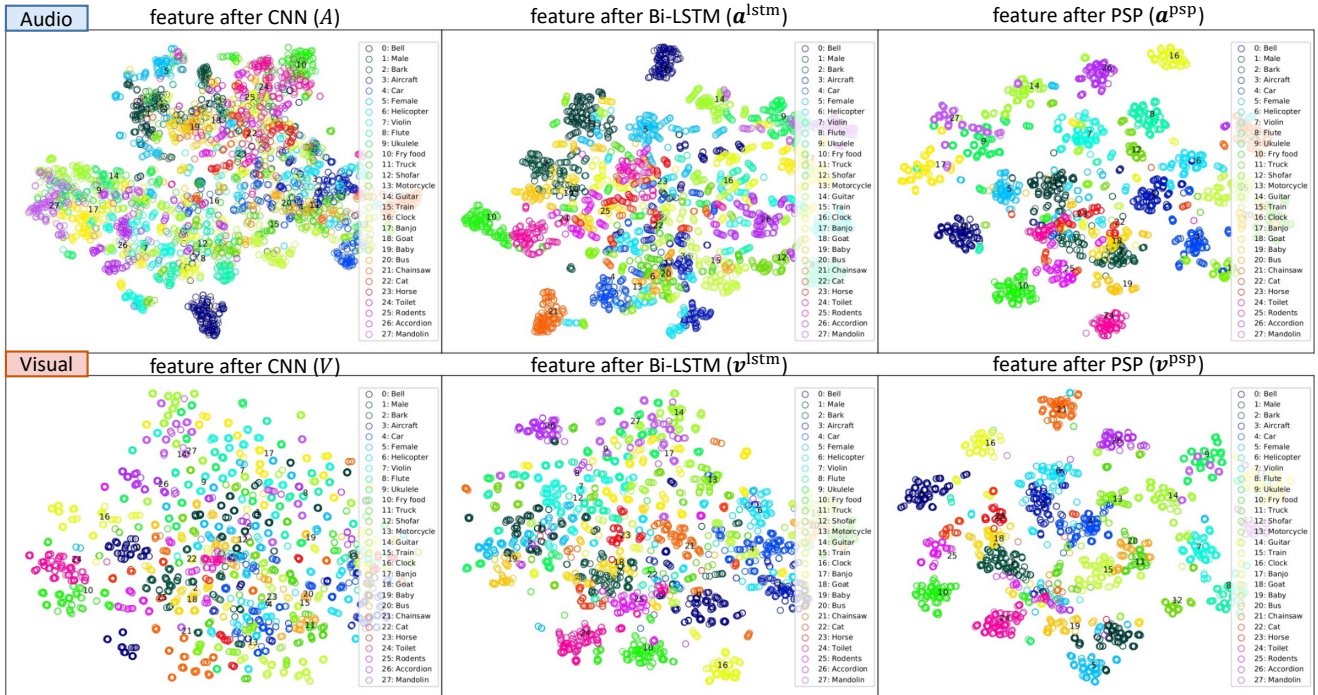


Figure 5. TSNE [16] visualization of audio and visual feature distributions under the fully-supervised setting. The data all come from the validation set. (Row 1:) audio features. (Row 2:) visual features. (Column 1:) the CNN features. (Column 2:) features after Bi-LSTM encoding. (Column 3:) features after PSP encoding. We observe that features after PSP are much better clustered into individual classes than the Bi-LSTM and CNN features. Different colors represent different classes. Best view in color and zoom in.

sions merely according to synchronized audio-visual segments while our method can pay attention to visual and audio components that are at different time stamps. For example, AVEL incorrectly regards the fifth and sixth segments as the *frying food* event, ignoring the third and fourth segments which are more relevant to the event. 3). We visualize the similarity matrices γ^{va} and γ^{av} in Fig. 4. We find that only a small percentage of all the audio-visual connections are retained after PSP selection and are closely related to the event. For example, for the first four visual components describing the target event, they tend to build strong connections (large similarity values) with the first three audio components containing the sound of the event. Such a propagation mechanism is critical for AVE localization because more discriminative audio-visual features can be identified with these *positive* connections and subsequently used in classifier training. Through backpropagation, it allows the model to be able to attend to broader and more sound-relevant regions in the visual images.

We then visualize the data distribution of features processed by different stages in our framework using TSNE [16] (Fig. 5). We first find that the CNN-based audio and visual features are not very well clustered. This is because they are at a relatively low level in the network hierarchy encoding limited semantics. Then, after Bi-LSTM, features of some categories (e.g., *rodents* and *Fry food*) can be better clustered compared with the CNN features, but

most are still disordered and highly mixed. Further, after PSP, the features are much better clustered: cohesive within the same class and divergent between different classes. This reflects that the audio-visual representations gain stronger discriminative abilities along the pipeline of our method.

6. Conclusion

For the AVE localization problem, we propose a positive sample propagation (PSP) method, which identifies and exploits relevant but unsynchronized audio and visual samples to enrich the encoded features. We find that negative and weak connections, even though having small weights, have a detrimental effect on the system, and thus have to be completely removed. Further, for the fully supervised setting, we propose an audio-visual pair similarity loss to supervise feature learning from a complementary way: whether a segment contains an event. For the weakly supervised setting, we insert a weighting branch to the classification module inject temporal importance to the features. Extensive experiments validate the effectiveness of these techniques.

Acknowledgement. This work was supported by the National Key Research and Development Program of China (2018YFB0804205), the National Natural Science Foundation of China (61725203, 61732008), the ARC Discovery Early Career Researcher Award (DE200101283), and the ARC Discovery Project (DP210102801).

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] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *ICCV*, 2017.
- [3] Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *ECCV*, 2018.
- [4] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *NeurIPS*, 2016.
- [5] Xuelian Cheng, Yiran Zhong, Yuchao Dai, Pan Ji, and Hongdong Li. Noise-aware unsupervised deep lidar-stereo fusion. In *CVPR*, 2019.
- [6] Ying Cheng, Ruize Wang, Zhihao Pan, Rui Feng, and Yuejie Zhang. Look, listen, and attend: Co-attention network for self-supervised audio-visual representation learning. In *ACM MM*, 2020.
- [7] Trevor Darrell, John W Fisher, and Paul Viola. Audio-visual segmentation and “the cocktail party effect”. In *ICMI*, 2000.
- [8] Haytham M. Fayek and Anurag Kumar. Large scale audiovisual learning of sounds with weakly labeled data. In *IJCAI*, 2020.
- [9] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *ICASSP*, 2017.
- [10] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *ICASSP*, 2017.
- [11] Di Hu, Feiping Nie, and Xuelong Li. Deep multimodal clustering for unsupervised audiovisual learning. In *CVPR*, 2019.
- [12] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv preprint arXiv:2102.05918*, 2021.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 2017.
- [14] Yan-Bo Lin, Yu-Jhe Li, and Yu-Chiang Frank Wang. Dual-modality seq2seq network for audio-visual event localization. In *ICASSP*, 2019.
- [15] Yan-Bo Lin and Yu-Chiang Frank Wang. Audiovisual transformer with instance attention for audio-visual event localization. In *ACCV*, 2020.
- [16] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 2008.
- [17] Sanjeev Parekh, Slim Essid, Alexey Ozerov, Ngoc QK Duong, Patrick Pérez, and Gaël Richard. Motion informed audio source separation. In *ICASSP*, 2017.
- [18] Jie Pu, Yannis Panagakis, Stavros Petridis, and Maja Pantic. Audio-visual object localization and separation using low-rank and sparsity. In *ICASSP*, 2017.
- [19] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, 2020.
- [20] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [21] Janani Ramaswamy. What makes the sound?: A dual-modality interacting network for audio-visual event localization. In *ICASSP*, 2020.
- [22] Janani Ramaswamy and Sukhendu Das. See the sound, hear the pixels. In *WACV*, 2020.
- [23] Mike Schuster and Kuldip K Paliwal. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, 1997.
- [24] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.
- [25] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In *CVPR*, 2018.
- [26] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [27] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-supervised audio-visual video parsing. In *ECCV*, 2020.
- [28] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *ECCV*, 2018.
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [30] Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention matching for audio-visual event localization. In *ICCV*, 2019.
- [31] Haoming Xu, Runhao Zeng, Qingyao Wu, Mingkui Tan, and Chuang Gan. Cross-modal relation-aware networks for audio-visual event localization. In *ACM MM*, 2020.
- [32] Hanyu Xuan, Zhenyu Zhang, Shuo Chen, Jian Yang, and Yan Yan. Cross-modal attention network for temporal inconsistent audio-visual event localization. In *AAAI*, 2020.
- [33] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018.