

Contrastive Positive Sample Propagation along the Audio-Visual Event Line

Jinxing Zhou, Dan Guo*, and Meng Wang*, *Fellow, IEEE*

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. It is pivotal to learn the discriminative features for each video segment. Unlike existing work focusing on audio-visual feature fusion, in this paper, we propose a new contrastive positive sample propagation (CPSP) method for better deep feature representation learning. The contribution of CPSP is to introduce the available full or weak label as a prior that constructs the exact positive-negative samples for contrastive learning. Specifically, the CPSP involves comprehensive contrastive constraints: pair-level positive sample propagation (PSP), segment-level and video-level positive sample activation (PSA_S and PSA_V). Three new contrastive objectives are proposed (*i.e.*, \mathcal{L}_{avpSP} , \mathcal{L}_{spSA} , and \mathcal{L}_{vpSA}) and introduced into both the fully and weakly supervised AVE localization. To draw a complete picture of the contrastive learning in AVE localization, we also study the self-supervised positive sample propagation (SSPSP). As a result, CPSP is more helpful to obtain the refined audio-visual features that are distinguishable from the negatives, thus benefiting the classifier prediction. Extensive experiments on the AVE and the newly collected VGGSound-AVEL100k datasets verify the effectiveness and generalization ability of our method.

Index Terms—Audio-visual event, positive sample propagation, contrastive learning, audio-visual learning.

1 INTRODUCTION

AN audio-visual event (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 sound it makes also exists in the audio portion (*audible*). The AVE localization task is to find these video segments which contain an audio-visual event and classify it into a certain category. We illustrate this task in Fig. 1. It belongs to the research topic of audio-visual scene understanding. It uses both audio and vision inputs to answer *if an event happens in both modalities at different video segments*. The task must explore unconstrained videos (events in real life) that are not limited to the temporal consistency of lip reading or other human-making sounds.

Recent literature has shown that by fusing multi-modality information can lead to better deep feature representation, *i.e.*, audio-visual fusion [1] and text-visual fusion [2]. However, building a large scale multi-modality pre-training datasets would require heavy manual labors to clean and annotate the raw video sets. To relief the manual labor, recent work either focuses on learning from noise supervision [3], [4] or tries to automatically filter out irrelevant samples [5]. In this work, we devote to explore better deep feature representation for AVEs, which is served for the latter purpose.

- J. Zhou, D. Guo, and M. Wang are with Key Laboratory of Knowledge Engineering with Big Data (HFUT), Ministry of Education, School of Computer Science and Information Engineering (School of Artificial Intelligence), Hefei University of Technology (HFUT), and Intelligent Interconnected Systems Laboratory of Anhui Province (HFUT), Hefei, 230601, China (e-mail: zhoujx@hfut.edu.cn; guodan@hfut.edu.cn; eric.mengwang@gmail.com).
- D. Guo and M. Wang are also with the Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei, 230601, China.
- Corresponding authors*: Dan Guo, Meng Wang.
- Code and dataset are available at <https://github.com/jasongief/CPSP>.

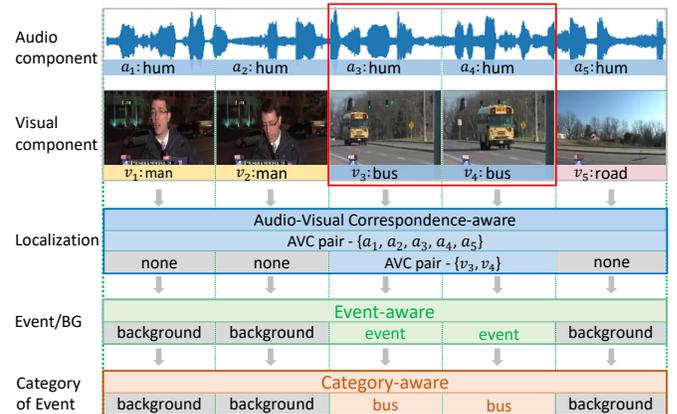


Fig. 1. An illustration of the AVE localization task. Each video segment is composed of an audio and a visual component. In this example, the “hum” of the bus exists in all the segments (audio modality), but the visual images of the “bus” appear only in the third and fourth segments (visual modality). So only these two segments (red boxes) are localized as *bus event*, the remaining are recognized as *background*. A localization system is expected to analyze and utilize the audio-visual pairs (audio-visual correspondence-aware), determine whether a video segment contains an audio-visual *event* (event-aware), and further identify its event category (category-aware).

We study how to effectively leverage audio and visual information for event localization. In the current AVE localization work, two relations are considered: intra-modal 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. Pioneer works [5], [6] often try to regress the class by directly concatenating features from synchronized audio-visual pairs; their accuracy is often unsatisfying. The following works [7]–[9] utilize a self-attention mechanism to explicitly encode the temporal relations within intra-modality and some of them [8]–[11]

also aggregate better audio-visual feature representations by encoding cross-modal relations. However, these methods aggregate all the audio and visual components, often ignore the interference caused by irrelevant audio-visual segment pairs during the fusion process. In this paper, we solve this problem in the cross-modal interaction by considering the relations from different perspectives: intra- and inter-videos. More importantly, we devote to feature aggregation and enhancement from high-relevant (positive) samples and can obtain better AVE localization accuracy.

About our observations, we argue and detail that the AVE localization task has three main challenges below. (1) *Unconstrained audio-visual relevance matching*. On one hand, the sound-maker is often occluded by some event-irrelevant objects or even be out of the screen, *e.g.*, the humming sound but accompanied by an announcer as shown in Fig. 1. On the other hand, there are usually multiple objects contained in the visual scene which could be sound-makers or not, *e.g.*, the humming accompanied by the bus, man, and road, also as shown in Fig. 1; the audio signal is also inevitably mixed with other noises. Such scenarios in unconstrained videos make it hard to match the audio-visual segment pairs in a flexible and accurate manner. 2) *Temporal inconsistency in AVE videos*. In real-life videos, the audio and visual signals are processed by an independent workflow. This spawns the research on the audio-visual synchronization problem [12]–[14] and brings the temporal inconsistency issue. Such issue lets the segment event judgment (*i.e.*, AVE or background) difficult. Especially for the weakly-supervised setting (refer Sec. 3 for the setting details), AVE localization task still asks to parse from segment-level but only given video-level label. (3) *Distinction of similar but different representations*. For the AVE localization task, it not only requires locating the event along the timeline but also must identify its category. In order to obtain discriminative feature representations, we must constraint the representations learning to be category-aware, such as distinguishing the videos displaying musical instruments guitar and violin although these two events are with a negligible difference.

The proposed Contrastive Positive Sample Propagation model (CPSP). To deal with aforementioned challenges, as shown in Fig. 2, we propose a new *CPSP method* that enables the localization system to encode discriminative representations by activating the *positive instances in the audio-visual data* from three levels, *i.e.*, the most relevant audio-visual pairs (*pair-level*), the segments containing an AVE consistently in audio and vision modalities (rather than the background, *segment-level*), the videos belonging to the same event category (rather than other categories, *video-level*). By exploiting these positive samples, the learned audio-visual representation encourages our model to be AVC (audio-visual correspondence)-aware, event-aware, and category-aware, which exactly matches the goal of AVE localization task. We introduce the details next.

Specifically, as shown in Fig. 3, we propose a new *Positive Sample Propagation (PSP) module*. In a nutshell, PSP 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,

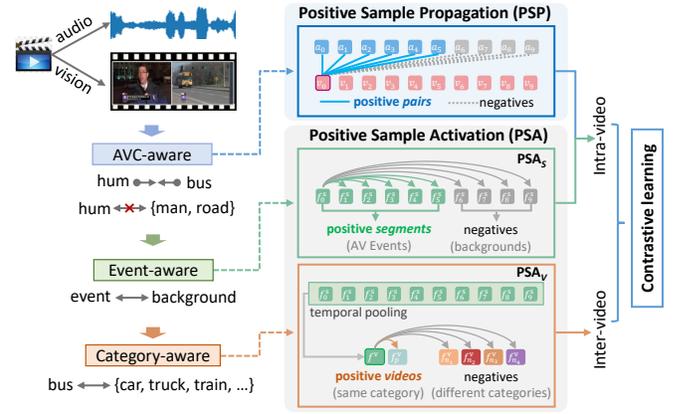


Fig. 2. An illustration of the proposed contrastive positive sample propagation (CPSP) method. The rounded rectangle in gray color represents a background segment, otherwise it means that the segment describes an event. The CPSP considers the AVE localization task from three aspects and works in a contrastive manner: 1) The *pair-level* positive sample propagation (PSP) aims to select the most relevant audio-visual segment pairs for cross-modal feature aggregation (*AVC-aware*). 2) The *segment-level* positive sample activation (PSA_S) aims to distinguish the positive segments that contain an audio-visual event from the background segments (*event-aware*). 3) The *video-level* positive sample activation (PSA_V) activates the videos sharing the same event category as the positive samples (*category-aware*).

we show that the PSP allows the most relevant features that are not necessarily synchronized to be aggregated in an online fashion. It is noteworthy that our PSP is *AVC-aware*, which is different from existing cross-modal feature fusion methods [8]–[11]. A concise illustration of the PSP has been shown in Fig. 2 and details can be seen from Fig. 4.

Apart from the PSP, it is also significant to effectively address the difficulties: 1) complex visual or audio backgrounds in an unconstrained video make it difficult to localize an AVE along the time line, and 2) localizing and recognizing an AVE category requires the model to deeply exploit the representative features among different AVE videos. Thus, as illustrated in Fig. 3, we propose a new *Positive Sample Activation (PSA) module* to constraint the representation of target segments containing an AVE to be possibly distinguishable from both the backgrounds in the same video (*intra-video*) and the segments of other videos with different event categories (*inter-video*). Specifically, the PSA is conducted from the segment-level (PSA_S) and video-level (PSA_V) whereas the PSA_S is designed to address the former purpose, while the PSA_V is proposed to solve the latter. As illustrated in Fig. 2, we show that the PSA_S is *event-aware* and the PSA_V is *category-aware*. Details are introduced in Sec. 4.3.

It is nontrivial to build positive and negative connections between complex visual scenes and intricate sounds. We utilize the principle of the *contrastive learning* without additional annotations to optimize the audio-visual representation learning. Three new contrastive losses, *i.e.*, the \mathcal{L}_{avpsp} for the PSP, \mathcal{L}_{spsa} for the PSA_S , and \mathcal{L}_{vpsa} for the PSA_V , are proposed for gathering positive samples while pushing the negatives away from them in the feature space. They are different from these audio-visual related works [15]–[17] that take the synchronized audio-visual segments as positive samples, which is not compatible with AVE localization. In our work, we explore the positive (high relevant) AVC

pairs and consider performing contrastive techniques from more levels, *i.e.*, adding the segment-level and video-level contrasting. To the best of our knowledge, we are the first to introduce contrastive learning to solve the AVE localization problem and provide a comprehensive discussion about the three levels of positive instances covering both intra- and inter- AVE video correlations as shown in fig. 2.

The improvement of backbone for fully and weakly supervised AVE localization. Besides the proposed CPSP, we analyze the fully and weakly supervised AVE localization, and further propose two improvements that work under each setting, respectively. On the one hand, an audio-visual pair similarity loss based on the PSP \mathcal{L}_{avpsp} is introduced under the fully supervised setting that encourages the network to learn high correlated features for the audio-visual pair if they contain 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 these two improvements on the standard AVE dataset [5] and the experimental results demonstrate their effectiveness. It is noteworthy that these two improvements are not only helpful in the proposed CPSP method but also can be applied to other localization networks (relevant results have been shown in Sec. 6.3.3).

To summarize, the proposed CPSP including the PSP and PSA modules devotes to better deep representation by imposing contrastive constraints on the localization system from three semantic levels (positive audio-visual instances of AVC pairs, segments, and videos), which covers both intra- and inter- AVE video correlations. The main contributions are summarized as follows:

- The proposed PSP allows to explore the most relevant (positive) *pair-level* features that are not necessarily synchronized but semantic-related to be aggregated and enables to encode more distinguishable audio-visual representations.
- The proposed PSA explicitly activates the positive samples from additional *segment-level* and *video-level* rather than directly sending the fused audio-visual features into the final classifier such as in previous works [8]–[11].
- The improvements of backbone proposed for the fully and weakly supervised settings consistently benefit our localization system and are also advantageous in other networks.
- Extensive experiments demonstrate the effectiveness of following designs and our method achieves new state-of-the-art performances under both settings.

At last, we remind that the PSP module was first introduced in our previous work [18]. Compared to the preliminary version, in this paper, we have made improvements in six aspects: (1) in addition to the PSP, we expand it to the CPSP by adding the Positive Sample Activation (PSA, including PSA_S and PSA_V) scheme that systematically exploiting segment-level and video-level positive audio-visual instances; (2) we perform a comprehensive survey of relevant works about the contrastive learning in the field of audio-visual representation learning in Sec. 2; (3) two new contrastive objective losses are designed and introduced into the AVE localization in Sec. 4.3; we add the discussion

about the objectives in Sec. 5; (4) we implement a self-supervised contrastive method SSPSP and give more analysis by comparing it with the CPSP in Sec. 6.4. (5) we release a large-scale VGGSound-AVEL100k dataset for AVEL task. The videos are sampled from VGGSound [19] where the video-level categories are given and we provide segment-level annotations. Considerable experiments are conducted on this large dataset to evaluate our model and more detailed analyses are provided in Sec. 6.4. (6) we extend the proposed method on the LLP [20] dataset collected for a similar but more challenging audio-visual video parsing (AVVP) task in Sec. 6.5 and the results also demonstrate the effectiveness and generalization ability of our method. In brief, the CPSP proposed in this paper makes the localization framework much more comprehensive and extensive experiments make the CPSP more convincing.

The rest of this paper is organized as follows. Sec. 2 provides an overview of the related works. Sec. 3 introduces two settings of the AVE localization problem, *i.e.*, the fully and weakly supervised tasks. Sec. 4 elaborates on the proposed CPSP. Discussions on the CPSP methodology are shown in Sec. 5. The experimental results and analyses are presented in Sec. 6, and conclusion are given in Sec. 7.

2 RELATED WORK

Audio-visual correspondence (AVC) aims to predict whether a given visual image corresponds to a short audio recording. The task 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 in video data. Existing methods try to evaluate the correspondences by measuring the audio-visual similarity [1], [21]–[24]. It will get a large similarity score if the audio-visual pair is corresponding, otherwise, a low score. This is in line with our focused AVE localization problem since the synchronized audio-visual pair of a target *event* segment must be corresponding. The difference is that AVC tackles the correspondence of an audio and an image, rather than an audio and a video in our AVE localization task. So, we are motivated to tackle the abundant audio-visual segment pairs in AVE localization problem by further exploring and exploiting the audio-visual similarity.

Sound source localization (SSL) aims to localize those visual regions which are relevant to the provided audio signal. It is related to *sound source separation* problem, which mainly focuses on the event of people speech [15], [25]–[29] or musical instrument playing [30]–[35]. For SSL, there is usually a condition that the sound source must appear in the visual image. In other words, it mainly focuses on the *visual* localization while the AVE localization devotes to the *temporal* localization. SSL has two settings: the single and multiple sound source(s) localization. For the single SSL, the localization map can be easier obtained in an unsupervised manner by directly computing the similarity between the audio feature and visual feature map. The multiple SSL is more challenging that requires to accurately locate the sound-maker when there are multiple sound sources [36]–[39]. The class activation mapping (CAM) [40] is helped to realize class-aware object localization. Qian *et al.* [36] adapt the Grad-CAM [41] to disentangle class-specific features for

multiple SSL. Hu *et al.* [37] maps audio and visual features into respective K cluster centers and take the center distance as a supervision to rank the paired audio-visual objects. Hu *et al.* [38] first learn the object semantics in single SSL then use that to help with multiple SSL. Recent methods use contrastive techniques to utilize the discriminative sound characteristics and diverse object appearances. Senocak *et al.* [42] propose a triplet loss working in an unsupervised manner. Afouras *et al.* [15] utilize a contrastive loss to train the model in a self-supervised learning way. Both of these methods [15], [42] need to construct positive and negative audio-visual pair samples. Considering the positive and negative samples can also be obtained in AVE localization task, we propose to exploit the abundant audio-visual instances for contrastive learning.

Audio-visual event localization (AVEL) aims to distinguish those segments including an audio-visual event from a long video. Different from the acoustic event classification [43]–[46] or video classification [47]–[51] making a prediction based on the whole audio or video embedding, AVE localization requires to judge the audio-visual correspondence and event category for each segment. Currently, the AVEL task is a supervised problem with weak labels or full labels. The former merely contains video-level labels, and the latter refers to both segment-level and video-level annotations. Existing works mainly focus on the audio-visual fusion process. A dual multimodal residual network is proposed in [5]. Lin *et al.* [6] adapt a bi-directional LSTM [52] to fuse audio and visual features in a seq2seq manner. These methods simply concatenate the synchronized audio-visual features during fusion. Subsequent work utilizes a bilinear method [10], [11] or a joint co-attention strategy [9], [53] to capture cross-modal relations between both synchronized and unsynchronized audio-visual pairs. Self-attention mechanism is also widely used to encode temporal relation in both the AVEL [7]–[9] and *audio-visual video parsing (a new weakly supervised audio-visual related task)* [20], [54], [55]. Lin *et al.* [56] 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.* [8] attempt to use the concatenating audio-visual features as the supervision then the feature of each modality is updated by separate modules. Unlike these, the proposed CPSP method has a further in-depth study on the abundant audio-visual pairs, event and background segments, and similar videos but with different categories, activating the most relevant ones. Relying on these positive samples, more distinguished audio-visual features can be obtained after feature aggregation.

Contrastive learning in audio-visual field. The technique of contrastive learning turns out to be an effective solution that is widely used in self-supervised learning [57]–[61] and various weakly supervised tasks [62]–[65]. Seeing data in a large batch size, the model learns discriminative representations by identifying the positive or negative samples. The key of contrastive learning is how to construct the positive and negative samples. Recently, some researchers start to explore injecting the label information to accurately select more positive/negative samples for better representation learning. Such supervised contrastive learning has

shown its superiority in both computer vision [66], [67] and natural language processing tasks [68], [69]. In the audio-visual related works [15]–[17], they usually adopt the self-supervised contrastive manner, *i.e.*, directly selecting the positive sample from the synchronized audio-visual segment and contrast them with the negatives come from different timestamps. However, this is not compatible with AVE localization: an audio and visual segment can be regarded as a positive sample as long as they describe the *same event*, and vice versa. To effectively distinguish an AVE, it is vital to construct exact positive and negative samples from the video segments for contrastive learning. Inspired by these, we propose the PSA in a supervised manner that explicitly performs contrastive learning between segments and videos with the segment/video-level label prior. To the best of our knowledge, we are the first to utilize the contrastive learning to solve AVE localization problem and we explore comprehensive contrastive strategies from different levels.

3 PROBLEM STATEMENT

AVE localization aims to find out those segments containing an audio-visual event [5]. 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 in length. 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 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 following [6], [9], 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 video segments that contain an event. This setting is different from the fully supervised one because the event label of each segment \mathbf{y}_t is unknown, making the problem more challenging.

4 OUR METHOD

The overall pipeline of our system is illustrated in Fig. 3, which includes four modules: a feature extraction and encoding module (Sec. 4.1), a positive sample propagation module (Sec. 4.2), a positive sample activation module (Sec. 4.3), and a classification module (Sec. 4.4.1). In the *feature extraction and encoding* module, the audio-guided visual attention (AVGA [5]) is adapted for early fusion to

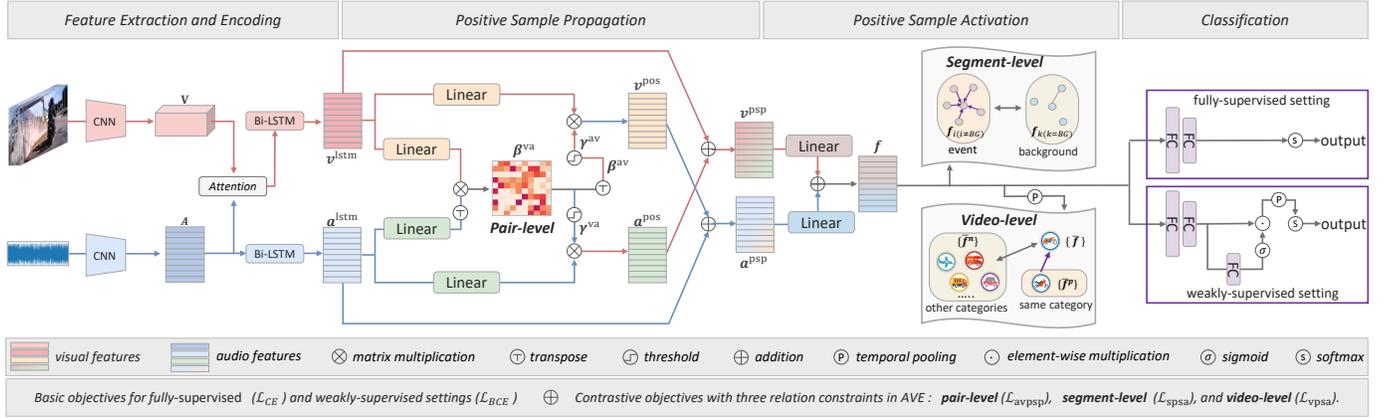


Fig. 3. System Flow. We first extract and encode video and audio features through existing modules such as AVGA [5] and Bi-LSTM. The proposed *positive sample propagation* (PSP) takes the LSTM encoded features as input; an affinity matrix is computed before selecting the connections of audio-visual segment pairs using thresholding. In this module, audio and visual features are aggregated by feature propagation through the *pair-level positive* connections. Then, the fused features are further processed under the constraints of the proposed *positive sample activation* (PSA), where the *segment-level* (PSA_S) enforces features of these video segments containing the same event (*positive*) to be possibly closer and be away from the background segments while the *video-level* (PSA_V) encourages representations of the videos sharing the same event category (*positive*) to be similar. Details are described in Sec. 4.3. In the last stage, we classify the event into predefined categories. For the supervised setting, apart from the commonly used cross-entropy (CE) loss, we further propose an audio-visual pair similarity loss based on the PSP 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. The whole system is optimized by the basic classification loss with additional proposed contrastive objectives, detailed in Sec. 4.4.2.

make the model focus on those visual regions closely related to the audio component. Then Bi-LSTM is utilized to encode temporal relations in video segments for each modality. The LSTM encoded audio and visual 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, *i.e.*, the pair-level contrastive constraint. Audio and visual features are aggregated by feature propagation through the positive connections. The fused audio-visual features after PSP are further processed by two contrastive constraints, *i.e.*, the *positive sample activation* from both segment-level (PSA_S) and video-level (PSA_V), which refine the audio-visual features for segments containing an event in a video and different videos but sharing the same event category, respectively. The updated features are then sent to the final *classification* module, predicting which video segments contain an event and the event category.

4.1 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 [5] 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.

4.2 Positive sample propagation (PSP)

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

In *all-pair connection construction*, all the audio-visual pairs are connected. As shown in Fig. 4, 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,

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

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. 4. 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,

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

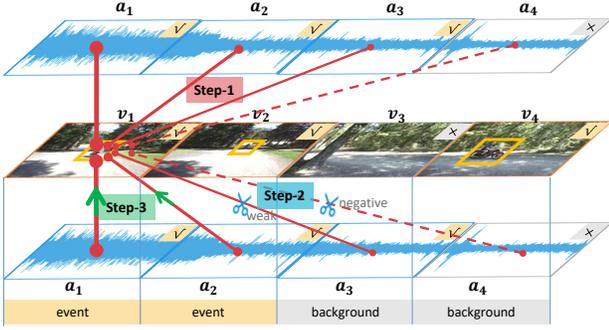


Fig. 4. An illustration of the proposed PSP. In this example, only the first two video segments contain an audio-visual event, *i.e.*, *motorcycle*. “√” denotes the audio or visual segment that describes the event, while “×” means not. The red lines denote connections of audio-visual pairs, solid lines 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. “Step-1” corresponds to the all-pair connection construction, while “Step-2” denotes the pruning of the negative and weak connections, and the green arrow indicates the *positive* direction of feature propagation in “Step-3”.

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. 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. 4. 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,

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

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

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. 9. More discussions are provided in Sec. 5.

4.3 Positive sample activation (PSA)

PSA is designed to make the model more event-aware and category-aware. It involves two steps. We activate more positive samples from both segment and video levels. We introduce the PSA_S and PSA_V with contrastive strategies below. Before that, we first fuse the audio and visual features \mathbf{a}^{PSP} and \mathbf{v}^{PSP} into an integrated audio-visual feature \mathbf{f} as follows:

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

where $\mathbf{f} \in \mathbb{R}^{T \times d_i}$ is the feature of video segments, $\mathcal{N}(\cdot)$ represents layer normalization, $\mathbf{W}_3^v, \mathbf{W}_3^a \in \mathbb{R}^{d_i \times d_i}$ represent learnable parameters in the linear layers. \mathbf{f} can be used to represent the segment feature and be summarized to the video feature.

4.3.1 Segment-level positive sample activation (PSA_S).

To make the model be event-aware, we design a contrastive strategy from the segment-level. As shown in Fig. 5(a), there are two sets of segments: segments depicting an audio-visual event constitute the *event set*, while remaining segments form the *background set*. We present a contrastive strategy to perceive the difference between these two video segment sets. Take arbitrary segment from the event set as an anchor, the remaining ones in the event set are regarded as its *positive* samples, while the segments in the background set are treated as *negative* samples. As shown in the right of Fig. 5(a), positive samples should be pulled together to the anchor, while the negative ones are pushed away. The segment-level contrastive objective takes the following form,

$$\begin{aligned} \mathcal{L}_{\text{spsa}} &= \\ &= -\frac{1}{N_i^e} \sum_{i=1}^{N_i^e} \log \left(\frac{\exp(\frac{\text{sim}(\mathbf{f}_i, \mathbf{f}_j)}{\eta})}{\exp(\frac{\text{sim}(\mathbf{f}_i, \mathbf{f}_j)}{\eta}) + \frac{1}{N_i^{\text{bg}}} \sum_{k=1}^{N_i^{\text{bg}}} \exp(\frac{\text{sim}(\mathbf{f}_i, \mathbf{f}_k)}{\eta})} \right), \end{aligned} \quad (5)$$

where features \mathbf{f}_i and \mathbf{f}_j belongs to the event set ($i \neq j$), \mathbf{f}_i is the segment anchor, \mathbf{f}_j is one of \mathbf{f}_i 's *positive* samples, \mathbf{f}_k denotes a *negative* sample comes from the background set. N_i^e and N_i^{bg} are the total numbers of the event and background segments, respectively. $\text{sim}(\cdot, \cdot)$ computes the dot product of the ℓ_2 normalized vectors (*i.e.*, cosine similarity); η is a temperature parameter controlling the concentration level of feature distribution.

4.3.2 Video-level positive sample activation (PSA_V).

To make the model be category-aware, we design another online contrastive strategy from video-level. As shown in Fig. 5(b), there are some instrument-related events in dataset (*e.g.* guitar, violin, mandolin, banjo, ukulele), which are hard to distinguish since they are similar in vision and sound. A main challenge for AVE localization is to correctly distinguish the event category.

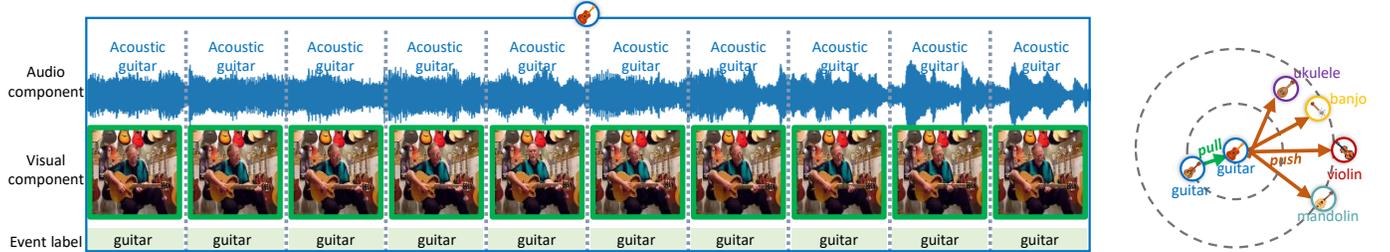
Specifically, for each data batch during training process, we compute the Euclidean distance between the video samples. Taking a video as an anchor, we select the *positive* sample identified with the same category and the largest distance. Similarly, we select *negative* videos that have different event categories from the anchor and the top- K closest distances, where K is a hyper-parameter controlling the number of *negative* samples. The model is expected to correctly recognize those similar but hard to learn samples: the positive sample should gather around the anchor video while the negative ones should be pushed farther.

To this end, the contrastive objective for video-level positive sample activation can be formulated as,

$$\mathcal{L}_{\text{vpsa}} = \max(0, d(\bar{\mathbf{f}}^a, \bar{\mathbf{f}}^p) - \frac{1}{K} \sum_{k=1}^K d(\bar{\mathbf{f}}^a, \bar{\mathbf{f}}_k^n) + \theta), \quad (6)$$



(a) Illustration of the *segment-level* positive sample activation (PSA_S). On the *left*, only the first five video segments contain the event *aircraft*; for these segments, take anyone as an *anchor*, the remaining four segments can be regarded as its *positive* samples, while the last five segments constitute the *negative* samples. Next, as illustrated on the *right*, in PSA_S , the positive samples are gathered around the anchor while pushed away from the negative ones.



(b) Illustration of the *video-level* positive sample activation (PSA_V). On the *left*, all the segments in the video contain the event *guitar*. Videos sharing the same event category *guitar* are treated as the *positive* samples, while the *negative* samples come from those holding other event categories. As illustrated on the *right*, for an anchor video, we merely select one positive sample belonging to the same event category with the largest Euclidean distance, and negative samples classified to other categories but with top- K smallest distances. The positive and top- K hard negative samples are online selected.

Fig. 5. We show two types of video data and illustrate the contrastive constraints (PSA_S and PSA_V) that are proposed to further exploit more positive samples from both segment-level and video-level, respectively. The green boxes represent an *event* happens in the visual segment, while gray boxes mean not.

where $d(\cdot, \cdot)$ computes the Euclidean distance between the ℓ_2 normalized vectors, $\bar{\mathbf{f}}^a$ is the feature vector of an anchor video, obtained by averaging feature of video segments \mathbf{f} (Eq. 4) along the temporal dimension, $\bar{\mathbf{f}}^p$ and $\bar{\mathbf{f}}^n$ are the features of positive and negative samples. K controls the number of the negative samples, θ denotes the minimum margin that the positive and negative samples should maintain.

4.4 Classification

The fused audio-visual feature $\mathbf{f} \in \mathbb{R}^{T \times d_i}$ is send to the classifier for prediction. We detail the classifier and the objective function for the fully and weakly supervised settings below.

4.4.1 Classifier

For the fully supervised setting, as shown in Fig. 3, the fused 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 [5], [6], [9], we add a weighting branch on the fully supervised classification module (Fig. 3). It is essentially another FC layer that enables the model to further capture the differences between synchronized audio-visual pairs. This process is summarized below,

$$\begin{cases} \mathbf{f}^h = \mathbf{f} \mathbf{W}_4^{\text{weak}} \mathbf{W}_5^{\text{weak}}, \\ \phi = \sigma(\mathbf{f}^h \mathbf{W}_6^{\text{weak}}), \\ \mathbf{o}^{\text{weak}} = s(\mathbf{f}_{\text{avg}}(\mathbf{f}^h \odot \Phi)), \end{cases} \quad (7)$$

where $\mathbf{W}_4^{\text{weak}} \in \mathbb{R}^{d_i \times d_h}$, $\mathbf{W}_5^{\text{weak}} \in \mathbb{R}^{d_h \times C}$, $\mathbf{W}_6^{\text{weak}} \in \mathbb{R}^{C \times 1}$ 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, \mathbf{f}_{avg} is the average operation along the temporal dimension. The final prediction $\mathbf{o}^{\text{weak}} \in \mathbb{R}^{1 \times C}$.

4.4.2 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 basic objective function, written as,

$$\mathcal{L}_{\text{ce}} = -\frac{1}{TC} \sum_{t=1}^T \sum_{c=1}^C \mathbf{Y}_{tc}^{\text{fully}} \log(\mathbf{O}_{tc}^{\text{fully}}). \quad (8)$$

Recall that each row of $\mathbf{Y}^{\text{fully}}$ contains a one-hot event label vector, describing the category of each video 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 based on the PSP $\mathcal{L}_{\text{avsp}}$. In principle, it asks the network to produce similar features for a pair of audio and visual components if the pair *contains an event* (contrasting from background) during PSP. 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

$$\mathbf{S} = \frac{\mathbf{v}^{\text{PSP}} \odot \mathbf{a}^{\text{PSP}}}{\|\mathbf{v}^{\text{PSP}} \odot \mathbf{a}^{\text{PSP}}\|_1}, \quad (9)$$

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

$$\mathcal{L}_{\text{avpsp}} = \mathcal{L}_{\text{MSE}}(\mathbf{S}, \mathbf{G}), \quad (10)$$

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

Combining Eq. 10 and Eq. 8, the objective function for fully-supervised setting $\mathcal{L}_{\text{fully}}$ can be computed by:

$$\mathcal{L}_{\text{fully}} = \mathcal{L}_{\text{ce}} + \lambda_1 \mathcal{L}_{\text{avpsp}}. \quad (11)$$

When refining the fused features with PSA_S and PSA_V jointly, the overall objective function $\mathcal{L}_{\text{fully}}^r$ can be computed by,

$$\mathcal{L}_{\text{fully}}^r = \mathcal{L}_{\text{fully}} + \lambda_2 \mathcal{L}_{\text{spsa}} + \lambda_3 \mathcal{L}_{\text{vpsa}}, \quad (12)$$

where $\lambda_1, \lambda_2, \lambda_3$ are hyper-parameters to balance the losses. The PSA_S and PSA_V can also be added to the vanilla PSP separately and such manner is slightly superior than joint training, we will discuss these two strategies in Sec. 6.4.1.

Weakly supervised setting. For this setting, following the practice in [6], [8], we adapt the binary cross entropy (BCE) loss as the basic classification loss, formulated as,

$$\mathcal{L}_{\text{weak}} = \mathcal{L}_{\text{BCE}}(\mathbf{o}^{\text{weak}}, \mathbf{Y}^{\text{weak}}). \quad (13)$$

It is worth mentioning that the segment-level event label should be known to distinguish the positive and negative segments in the video, thus, PSA_S is not applicable for weakly supervised AVE localization. The PSA_V can be used in both fully and weakly supervised settings since both of them provide the video-level event label. When combined with the loss item of PSA_V , the overall objective can be written as,

$$\mathcal{L}_{\text{weak}}^r = \mathcal{L}_{\text{weak}} + \lambda_4 \mathcal{L}_{\text{vpsa}}, \quad (14)$$

where λ_4 is a balance weight.

5 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 that \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. 9.

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 of 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 $\mathcal{L}_{\text{avpsp}}$. As shown in Eq. 8, the classification loss \mathcal{L}_{ce} prompts the model to calculate the loss between the output probabilities and the ground truth label. In comparison, $\mathcal{L}_{\text{avpsp}}$ allows the network to be aware of *whether an event exists in an audio-visual pair* (pair-level contrasting). 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, $\mathcal{L}_{\text{avpsp}}$ provides another auxiliary constraint so that the model can better select the most relevant visual (audio) components for feature aggregation during PSP. Note that $\mathcal{L}_{\text{avpsp}}$ cannot be adapted to the weakly supervised setting, where the label g_t of each segment is unknown. To summarize, \mathcal{L}_{ce} and $\mathcal{L}_{\text{avpsp}}$ serve as strong supervisions, especially in the fully supervised setting.

Implications of $\mathcal{L}_{\text{spsa}}$ and $\mathcal{L}_{\text{vpsa}}$. The design of $\mathcal{L}_{\text{spsa}}$ in Eq. 5 and $\mathcal{L}_{\text{vpsa}}$ in Eq. 6 are intended to further advance the audio-visual representation learning. $\mathcal{L}_{\text{spsa}}$ allows the network to be aware of *whether an event exists in a segment unit*, and $\mathcal{L}_{\text{vpsa}}$ allows the network to be aware of *whether an exact event category exists in a video*. Here, the network is enforced to learn the discriminative representation capability with the most relevant segments and the same category samples. In fact, these two losses can be regarded as the soft supervisions since they are controlled by the hyper-parameters (*i.e.*, η for $\mathcal{L}_{\text{spsa}}$, K and θ for $\mathcal{L}_{\text{vpsa}}$). The positive segment and video samples respectively activated by the PSA_S ($\mathcal{L}_{\text{spsa}}$) and PSA_V ($\mathcal{L}_{\text{vpsa}}$) are beneficial for the classifier training and this can be confirmed by the visualization examples shown in Figs. 10, 11. To summarize, $\mathcal{L}_{\text{avpsp}}$ (PSP) and $\mathcal{L}_{\text{spsa}}$ (PSA_S) exploit the positive clues of intra-video correlation (using pair and segment event labels), $\mathcal{L}_{\text{vpsa}}$ (PSA_V) focuses on the positive clues of inter-video correlation (using video category label). As the same reason for $\mathcal{L}_{\text{avpsp}}$, $\mathcal{L}_{\text{spsa}}$ is used for fully supervised setting while $\mathcal{L}_{\text{vpsa}}$ is not limited to this constraint.

6 EXPERIMENT

6.1 Experimental setup

Dataset. (1) AVE dataset [5]. Following the existing works [5], [6], [8], [9], we use the public AVE dataset for localization. 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 the duration of each segment is one-second. The audio-visual event boundary on the segment level and the event category on the video level are provided. Keeping consistent with prior work, 3,339 videos are used for training, while both the validation and test set contains 402 videos. **(2) VGGSound-AVEL100k dataset.** We construct a new large-scale VGGSound-AVEL100k dataset for AVEL task, in which the videos are sampled from VGGSound [19]. VGGSound-AVEL100k contains 101,072 videos that spans 141 audio-visual event categories covering more scenes in real-life that do not appear in the AVE dataset, such as motorboat, electric shaver, sharpen knife, *etc.* The ratio of train/validation/test split percentages are set as 60/20/20¹. **(3) LLP dataset [20].** It is collected for a more challenging audio-visual video parsing (AVVP) task where only video-level labels are given. It contains 11,849 videos and each video in this dataset contains multiple audio and visual events. The AVVP task requires to predict what events happen in both audio and visual tracks, separately. We extend the proposed CPSP in the weakly supervised setting to this task to evaluate the generalization ability of our method.

Evaluation metric. The category label of each segment is predicted in both fully and weakly supervised settings. Following [5], [6], [8], [9], we adopt the classification accuracy of each segment as the evaluation metric.

Training procedure and configuration. We have to deal with all the video data of different types (as shown in Fig. 5). For convenience, we use \mathcal{D}_{bg} to denote this type of video that contains both AVE and background segments (Fig. 5(a)) and use \mathcal{D}_{ae} to represent another type of video contains pure AVE segments belonging to a certain event category (Fig. 5(b)). We initialize the proposed localization system (Fig. 3) with the objective function \mathcal{L}_{fully} (Eq. 11) and \mathcal{L}_{weak} (Eq. 13) on the dataset benchmark (\mathcal{D}_{bg} & \mathcal{D}_{ae}). Then we further refine the audio-visual features with PSA_S on subset \mathcal{D}_{bg} and PSA_V on subset \mathcal{D}_{ae} in Eqs. 5 and 6, respectively; their corresponding usages (objective functions) for fully and weakly supervised settings are introduced in Eqs. 12 and 14. More related details are discussed in Sec. 6.3.2. We abbreviate the initialized network as **PSP**, and the refined network as **CPSP** in the following experiment evaluation. The parameters are tuned on the validation set with the final model is tested on a held-out test set. The results are reported in Sec. 6.3 and 6.4.

Implementation details. (1) *Visual feature extractor.* For fair comparison, we use the VGG-19 [70] pretrained on ImageNet [71] to extract the visual features. Specifically, 16 frames are sampled from each one-second video segment. We extract the visual feature map for each frame from the *pool-5* layer in VGG-19 with the size of $7 \times 7 \times 512$ and then use the average map as the visual feature for this segment. (2) *Audio feature extractor.* For audio features, we first process the raw audio into log-mel spectrograms and then use the VGGish, a VGG-like network [43] pretrained on AudioSet [72], to extract the acoustic feature with the dimension of 128. (3)

TABLE 1

Comparison with the state-of-the-art methods under both the fully and weakly supervised settings. We report the accuracy(%) measured on the AVE and the VGGSound-AVEL100k datasets. * indicates the number is reproduced by us.

Method	AVE		VGGSound-AVEL100k	
	fully	weakly	fully	weakly
AVEL [5]	68.6	66.7	55.7*	46.2*
AVSDN [6]	72.6*	67.3*	-	-
CMAN [9]	73.3*	70.4*	-	-
DAM [7]	74.5	-	-	-
AVRB [11]	74.8	68.9	-	-
AVIN [10]	75.2	69.4	-	-
RFJCA [53]	76.2	-	-	-
AVT [56]	76.8	70.2	-	-
CMRA [8]	77.4	72.9	57.1*	46.8*
MPN [74]	77.6	72.0	-	-
PSP [18](Ours)	77.8	73.5	58.3	47.4
CPSP(Ours)	78.6	74.2	59.9	48.4

Hyper-parameter settings. The temperature parameter η in Eq. 5 is set to 0.1. Impacts of the number of negative samples K and the margin θ in Eq. 6 are discussed in Sec. 6.3.2. The weights of λ_1 in Eq. 11 and λ_2, λ_3 in Eq. 12 is empirically set to 100, 0.01, 1, respectively. λ_4 in Eq. 14 is set to 0.005. These hyper-parameters remain the same on the AVE and the VGGSound-AVEL100k datasets in our experiments. In addition, the batch size in our experiments is set to 128. We use Adam [73] as the optimizer, and dropout technique is used in all the linear layers (Fig. 3) with the drop rate set to 0.1. As for the experiments on LLP dataset for video parsing, the batch size is set to 16 same as baseline method HAN [20], more implementation details are introduced in Sec. 6.5.

6.2 Comparison with state of the arts

We compare our method with the state of the arts in Table 1 by evaluating on the AVE and the VGGSound-AVEL100k datasets. Taking the results on the AVE dataset for example, compared with the baseline method AVEL [5], the CPSP exceeds it by 10.0% and 7.5% under the fully and weakly supervised settings, respectively. Such superiority is also proved by the results shown in Table 6. Also, our method exceeds those SOTAs [8]–[11] that focus on the cross-modal feature fusion using all of the audio-visual pairs. This indicates the necessity of the positive pair selection in PSP. CMRA [8] has comparable performance with the PSP method, but the CPSP is superior than CMRA on both datasets in both settings. Also, the CPSP exceeds the PSP method by a large margin. This again demonstrates the effectiveness of the PSA performing additional contrastive learning from both segment-level and video-level. Such advantages of the proposed CPSP method can also be observed on the large-scale VGGSound-AVEL100k dataset. For example, the CPSP exceeds the competitive CMRA and vanilla PSP by a large margin. We also notice that all the methods have a performance drop on the large-scale VGGSound-AVEL100k compared to the AVE dataset, we speculate VGGSound-AVEL100k contains much more videos with more event categories that makes the problem more challenging. Nevertheless, our method is more superior and

1. The large-scale VGGSound-AVEL100k dataset for AVEL task is available at <https://drive.google.com/drive/folders/1en1dks1GYiGaDS9ArQj1mmyoOdzEsQj?usp=sharing>. We give more details in Appendix. A

TABLE 2

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.

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 3

Impact of various values of τ on the system accuracy evaluated on the AVE dataset. Results on the two settings are shown.

τ	0	0.025	0.075	0.095	0.115
Fully-supervised	76.0	76.1	75.3	77.8	76.6
Weakly-supervised	71.2	71.7	70.4	73.5	72.8

robust under all the settings, which can be attributed to our system design.

6.3 Quantitative analysis - main modules

Here we test the effects of the PSP and PSA modules. Ablation experiments are mainly conducted on the AVE dataset.

6.3.1 Evaluation of the proposed PSP module

The effectiveness of the PSP encoding can be verified through an ablation study in Table 2. In Table 2, we denote the method without PSP, *i.e.*, removing it from the localization network (Fig. 3), as “w/o PSP”. We observe from the table that the performance on the AVE dataset drops in both the fully supervised and weakly supervised settings 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 the PSP.

Comparison with alternative pair-level 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 (denoted as “ASP”); (2) only negative ones are removed, while weak connections are remained (denoted as “WPSP”). Results are shown in Table 2. We have two main observations. First, when all samples are propagated, the accuracy of “ASP” drops by 1.9% and 2.3% on the fully and weakly supervised settings, respectively. This shows that it is essential to have a selection process before feature aggregation instead of utilizing all the connections. Second, although we merely remove the negative connections (*i.e.*, with a similarity value below $\tau = 0$), the system of “WPSP” is inferior to the full method. Specifically, the classification accuracy decreases by 1.8% and 2.3% under the fully and weakly supervised settings, which validates the effectiveness of filtering out the negative connections.

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 the accuracy generally 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. Such a connection pruning (*i.e.*, positive pair selection) process in PSP can be clearly observed from the visualization example in Fig. 8.

Comparison with adding self-attention [75] to the feature extractor. Self-attention [75] is widely used in existing methods [7]–[9], [20] 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 in the feature extractor module and denote it as the “SAPSP” method. As shown in Table 2, the performance surprisingly decreases by 2.4% and 2.7% under fully and weakly supervised settings, respectively. This indicates that in our system, it is not required to add additional intra-modal verification through self-attention before the PSP module. We speculate that the PSP is sufficient to describe the cross-modality while implicitly reveals the intra-modality correlations.

6.3.2 Evaluation of the proposed PSA module

Effectiveness of the PSA_S . PSA_S is expected to constraint the model to learn consistent features for the video segments containing the same event, while possibly be distinguishable from the background segments. We reflect its effect by the distance of centroids of the event and background segments in feature space. For videos in the dataset, we encode the segment features by the PSP and CPSP (merely equipped with PSA_S for fair comparison), respectively. Specifically, for each same event category, we filter out and average the features of event and background segments respectively; we take the two obtained vectors as event and background centroids. We calculate their Euclidean distance. Results on the AVE dataset are presented in Fig. 6. We can see that the distances between event and background segments are increased in most of the categories (24 out of 28) using the CPSP method. For example, for the event of *female, guitar*, and *bus*, the distances increase by around 33%. This verifies the benefit of the PSA_S that activates the event segments such that they can be better recognized from the backgrounds.

Effectiveness of the PSA_V . As introduced in Sec. 4.3, PSA_V aims to distinguish the positive video from the top- K closest but negative samples, and the Euclidean distance between the video-level representations is expected to be no less than the margin θ (Eq. 6). Here, we test CPSP merely with PSA_V for fair comparison. We conduct a study on the AVE dataset to explore the impacts of parameters K and θ in PSA_V . First, we empirically fix the number of negatives K to 4 and sample θ from {0.2, 0.4, 0.6}. As shown in Table 4, the performance is gradually improved as θ increases. And the best accuracy is achieved when $\theta = 0.6$ for both settings (*i.e.*, 78.31% for fully supervised, 74.20% for weakly supervised). We speculate that this is a relatively large margin to better distinguish the positive and negative samples. Next, we fix θ to 0.6 and test K with values {1, 2, 4, 6}. As observed from the table, $K = 4$ is the optimal setup. This means four negative video samples are selected from a batch of data

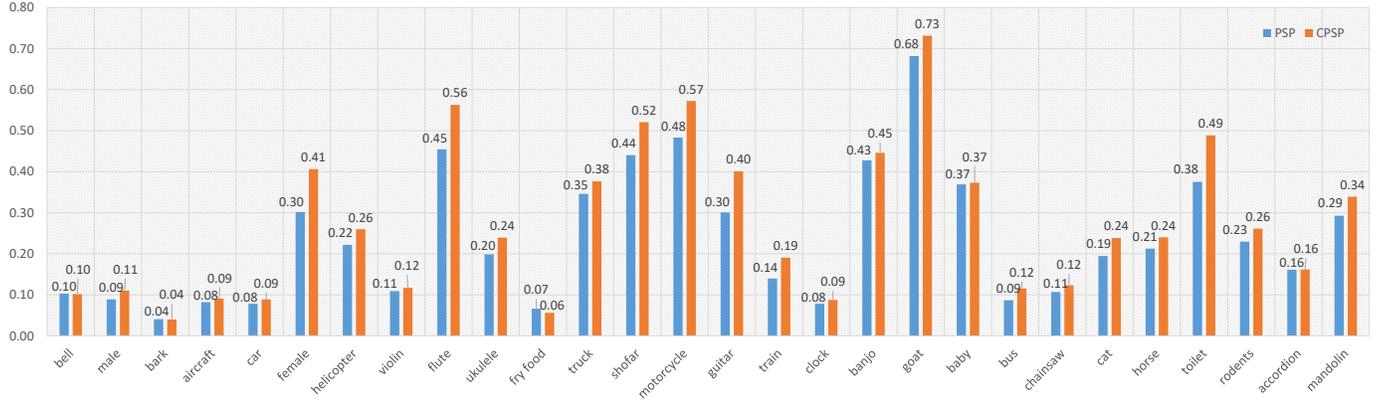


Fig. 6. Euclidean distances between the centroids of event and the background segments for each event category in the fully supervised setting. We respectively evaluate the segment features learned by the PSP and CPSP (merely w. PSA_S) in fully supervised setting. Larger distance of the CPSP demonstrates the benefit of PSA_S helping to encode features of event and background that are easier to distinguish. This experiment is conducted on the AVE dataset.

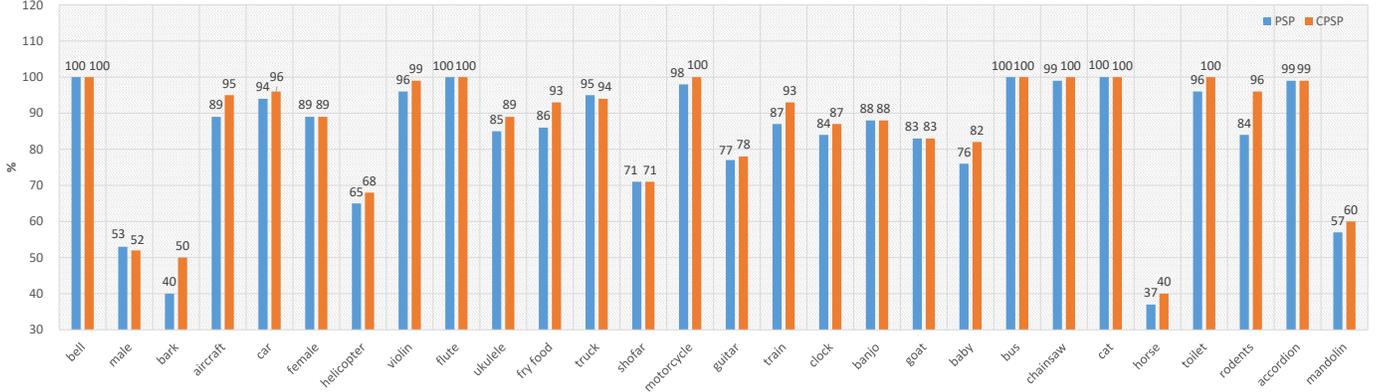


Fig. 7. Classification accuracy for videos containing no background segments in the fully supervised setting. Compared with the PSP, the CPSP (merely w. PSA_V) has overall superior performance in predicting the correct event categories. This experiment is conducted on the AVE dataset.

TABLE 4

Parameter study of the K and θ in PSA_V . We report the performance of the CPSP under different PSA_V setups in both fully and weakly supervised settings. Experiments are conducted on the AVE dataset. The **bold-faced** results represent the optimal performance is achieved under that setup.

Parameter setup		Fully-supervised	Weakly-supervised
$K = 4$	$\theta = 0.2$	77.61	74.10
	$\theta = 0.4$	78.10	74.18
	$\theta = 0.6$	78.31	74.20
$\theta = 0.6$	$K = 1$	78.20	74.15
	$K = 2$	78.13	74.15
	$K = 4$	78.31	74.20
	$K = 6$	78.23	74.10

during training to compare with the positive one. In this way, the model enables to simultaneously compare videos of multiple event categories at once in an online fashion. We set $K = 4$ and $\theta = 0.6$ for all of our other experiments on both the AVE and the VGGSound-AVEL100k dataset when conducting PSA_V . It is worthy to note that the performances of almost all the setups in the CPSP exceed the results of the PSP (*i.e.*, obtained from the case without PSA_V , 77.8% and 73.5% accuracy for fully and weakly supervised settings, respectively).

To clarify the effect of PSA_V more clearly, here we report the classification accuracy for each event category under the fully supervised setting on subset \mathcal{D}_{ae} of AVE dataset, where videos in \mathcal{D}_{ae} contain no background segments (*i.e.* having definite video-level category). The results are shown

in Fig. 7, the CPSP (merely equipped with PSA_V here) has better performance in most event categories (26 out of 28). This verifies that PSA_V can further help to predict the accurate event category, which is contributed to the video-level contrastive learning that makes the learned features more distinguishable for videos owing to different categories.

6.3.3 Evaluation of the improvement in fully/weakly setting Effectiveness of the pair similarity loss \mathcal{L}_{avpsp} in fully supervised setting. We respectively adapt \mathcal{L}_{ce} and $\mathcal{L}_{ce} + \lambda_1 \mathcal{L}_{avpsp}$ as the objective function and test them for model training. Two baselines are used: our PSP system and the AVEL system [5]. Results are presented in Table 6. We can clearly see that \mathcal{L}_{avpsp} improves the accuracy when the system is fully supervised. The improvement is 1.2% and 1.5% for PSP and AVEL, respectively. These results confirm the role of \mathcal{L}_{avpsp} 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 [5], [6], [9] consists in the weighting branch (Fig. 3). 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 6. We find that the performance of PSP and AVEL is improved by 1.9% and 2.3%, respectively. We argue that the additional weighting branch within the designed classification module allows the model to give different weights to the temporal sequences, thus benefiting

TABLE 5

Results on the AVE and VGGSound-AVEL100k datasets are reported. We list details of the experimental configurations, *i.e.*, the objective function (Objective), the type of video data used in the objective optimization (Data), the initialized model (Init.), the trained model (Return), and the learning rate (Lr).

Fully-supervised setting									
Method	Objective			Data	Init.	Return	Lr	Accuracy	
	$\mathcal{L}_{ce} \& \mathcal{L}_{avpsp}$	\mathcal{L}_{spsa}	\mathcal{L}_{vpsa}					AVE	VGGSound-AVEL100k
PSP	✓			$\mathcal{D}_{bg} \& \mathcal{D}_{ae}$	Xavier [76]	\mathcal{M}_{fully}^p	10^{-3}	77.8	58.3
CPSP _S	✓	✓		\mathcal{D}_{bg}	\mathcal{M}_{fully}^p	\mathcal{M}_{fully}^{sp}	10^{-4}	78.2	59.6
CPSP _V	✓		✓	\mathcal{D}_{ae}	\mathcal{M}_{fully}^p	\mathcal{M}_{fully}^{cp}	10^{-5}	78.3	59.8
CPSP(join)	✓	✓	✓	$\mathcal{D}_{bg} \& \mathcal{D}_{ae}$	\mathcal{M}_{fully}^p	$\mathcal{M}_{fully}^{osp}$	10^{-5}	78.4	59.8
CPSP(sepa)	✓	✓	✓	$\mathcal{D}_{bg} \rightarrow \mathcal{D}_{ae}$	\mathcal{M}_{fully}^{sp}	$\mathcal{M}_{fully}^{osp}$	10^{-5}	78.6	59.9
Weakly-supervised setting									
Method	Objective		Data	Init.	Return	Lr	Accuracy		
	\mathcal{L}_{bce}	\mathcal{L}_{vpsa}					AVE	VGGSound-AVEL100k	
PSP	✓		$\mathcal{D}_{bg} \& \mathcal{D}_{ae}$	Xavier [76]	\mathcal{M}_{weak}^p	10^{-3}	73.5	47.4	
CPSP	✓	✓	\mathcal{D}_{ae}	\mathcal{M}_{weak}^p	\mathcal{M}_{weak}^{cp}	10^{-5}	74.2	48.4	

TABLE 6

Method comparison on the AVE dataset under two settings. We evaluate 1) loss \mathcal{L}_{avpsp} 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 [5]. Under AVEL, * denotes that the number is produced by us. We use **bold** font to show the higher performance brought by our technique.

Setting	Method	PSP [18](ours)	AVEL [5]
fully	\mathcal{L}_{ce}	76.6	69.8*
	$\mathcal{L}_{ce} + \lambda_1 \mathcal{L}_{avpsp}$	77.8	71.3*
weakly	w/o weight. branch	71.6	66.9*
	w/ weight. branch	73.5	69.2*

the localization of the target video segments. These results confirm the effectiveness of the proposed improvements and show their robustness in other localization network. We refer readers to Sec. 5 for methodological discussions on the this technique.

6.4 Quantitative analysis - contrastive manner in CPSP

In this subsection, we first compare the PSP (without contrastive learning) and the CPSP (with contrastive learning), and then discuss the supervised CPSP and self-supervised PSP (named SSPSP).

6.4.1 Comparison of the PSP and CPSP

To reveal the impact of each contrastive loss, we list five training modes in Table 5: ① the vanilla PSP [18], ② the CPSP with sole PSA_S (denoted as **CPSP_S**), ③ the CPSP with sole PSA_V (denoted as **CPSP_V**), ④ the CPSP with both PSA_S and PSA_V by jointly training (denoted as **CPSP(join)**), and ⑤ the CPSP with both PSA_S and PSA_V by separately training (first PSA_S then PSA_V, denoted as **CPSP(sepa)**).

As shown in Table 5: (1) Compared with the vanilla PSP, the performances are improved by the CPSP with PSA_S and PSA_V in both fully and weakly supervised settings for both datasets. Take the VGGSound-AVEL100k dataset for example, after utilizing the PSA_S or PSA_V separately for the fully supervised setting, the accuracy of **CPSP_S** and **CPSP_V** increase by 1.3% and 1.5% (from 58.3% to 59.6% and 59.8%), respectively; (2) When jointing the PSA_S and

TABLE 7

Quality analysis of the encoded features of videos in the AVE dataset. We measure three clustering metrics (SC, CH, DBI) in two ways: "Mean" denotes evaluating two clusters (event and background) in each event category, while "All" refers to $C - 1$ clusters, where $C - 1$ is the number of all the event categories except background.

Method	SC ↑		CH ↑		DBI ↓	
	Mean	All	Mean	All	Mean	All
PSP	0.17	0.19	16.64	194.12	1.62	2.07
CPSP	0.21	0.22	18.29	195.51	1.54	1.95

PSA_V together, the performances keep stable under both combination strategies, *i.e.*, **CPSP(join)** and **CPSP(sepa)**. The results of **CPSP(join)** and **CPSP(sepa)** are comparable and **CPSP(join)** performs slightly worse. We argue that it may be a little confusing for **CPSP(join)** to train with the PSA_S and PSA_V simultaneously (totally different contrastive goals). With **CPSP(sepa)**, the best accuracy can achieve 78.6% and 59.9% for the AVE and the VGGSound-AVEL100k datasets, respectively. The CPSP is flexible to be utilized: with single independent PSA_S or PSA_V module, or with both modules under either training setup. Anyway, these benefit from the proposed PSA by activating the model to distinguish (1) positive and negative segments in PSA_S, and (2) event categories in PSA_V, the model gains a more robust capability to correctly identify the event location and its categories. This is consistent with the goal of AVE localization thus can promise better results, which can also be verified by the qualitative examples as shown in Figs. 10, 11.

Moreover, we introduce three widely-used clustering metrics, *i.e.*, Silhouette Coefficient (SC) [77], Calinski-Harabasz Index (CH) [78], and Davies-Bouldin Index (DBI) [79]. These metrics validate the data clustering quality from the intra-class aggregation and inter-class separation. Here, we use them to evaluate the event aggregation and background separation of features learned by the PSP and CPSP. At first, we have a brief introduction: (1) SC [77] calculates the difference of intra-class and inter-class dissimilarities and divides it by the maximum value of these two dissimilarities; larger score denotes that clusters are dense while well separated; (2) CH [78] is the ratio of the covariance of the intra-class

TABLE 8
Comparison with the baseline methods on the test set of LLP dataset. † denotes the results reported in the paper HAN [20].

Method	Segment-level					Event-level				
	A	V	AV	Type@AV	Event@AV	A	V	AV	Type@AV	Event@AV
AVEL [5] †	47.2	37.1	35.4	39.9	41.6	40.4	34.7	31.6	35.5	36.5
AVSDN [6] †	47.8	52.0	37.1	45.7	50.8	34.1	46.3	26.5	35.6	37.7
HAN [20]	60.1	52.9	48.9	54.0	55.4	51.3	48.9	43.0	47.7	48.0
PSP [18]	54.2	54.7	48.3	52.4	52.5	46.8	50.2	42.8	46.6	45.6
CPSP (Ours)	58.5	57.8	52.6	56.3	55.8	51.6	54.0	46.5	50.7	49.9

data to the covariance of the inter-class data; higher score means better performance; (3) DBI [79] represents the average similarity between clusters; a lower value indicates better separation between clusters. Next, we measure these metrics (*i.e.*, SC, CH, and DBI) in two ways: as shown in Table 7, the “Mean” denotes that we first split the data into event and background segments (*i.e.*, 2 clusters, binary separation) in each event category, and then average the metrics over all the categories; the “All” refers to $C - 1$ clusters, including all the different event categories except background (*i.e.*, multi-class event separation). In other words, we adopt “Mean” to measure the clustering effect with the partition of event and background, while “All” with the partition of all the event categories. At last, experiment is performed on the AVE dataset. As observed from Table 7, for all of the metrics, CPSP is better scored than PSP in any measurement method. This demonstrates that the positive features learned through CPSP are better clustered thus making it easier to distinguish the event segments from backgrounds, and also performing better for classifying videos with different categories.

6.4.2 Comparison of the CPSP and SSPSP

The contrastive learning is always conducted in a self-supervised manner in the audio-visual field [15]–[17]. Specifically, the synchronized audio-visual segment pair is regarded as a positive sample and otherwise is negative. There is a drawback that this manner will inevitably bring false negatives of the audio-visual pair depicting the same event (existing AVC) but at different timestamp. We are curious about the effect of using such self-supervised manner in AVE localization task. So we introduce the self-supervised learning into our positive sample propagation, and denote it as “SSPSP” method. In “SSPSP”, all unsynchronized features (*i.e.*, audio feature \mathbf{a}^{PSP} and visual feature \mathbf{v}^{PSP}) are treated as negative instances sampling from a batch of data during training. The corresponding contrastive objective can be written as \mathcal{L}_{ss} below. We first train a vanilla PSP with $\mathcal{L}_{\text{fully}}$ (Eq. 11), and then inject the self-supervised learning to the PSP. The total objective function can be computed by

$$\begin{cases} \mathcal{L}_{\text{sspsp}} = \mathcal{L}_{\text{ce}} + \lambda'_2 \mathcal{L}_{\text{ss}}, \\ \mathcal{L}_{\text{ss}} = -\frac{1}{N^b T} \sum_{i=1}^{N^b T} \log\left(\frac{\exp(\frac{\text{sim}(\mathbf{a}_i^{\text{PSP}}, \mathbf{v}_i^{\text{PSP}})}{\eta})}{\sum_{j=1}^{N^b T} \exp(\frac{\text{sim}(\mathbf{a}_i^{\text{PSP}}, \mathbf{v}_j^{\text{PSP}})}{\eta})}\right), \end{cases} \quad (15)$$

where N^b denotes the number of videos in a batch and T is the number of segment in each video; thus, $N^b \cdot T$ denotes the total segment number in a batch. $\text{sim}(\cdot, \cdot)$ computes the dot product of the ℓ_2 normalized vectors (*i.e.*, cosine similarity). The i and j are the indexes of the video segment, note that j

TABLE 9
Comparison of different contrastive manners (CPSP v.s. self-supervised SSPSP) under two settings. We report the accuracy(%) measured on the AVE and the VGGSound-AVEL100k datasets.

Method	AVE		VGGSound-AVEL100k	
	fully	weakly	fully	weakly
PSP	77.8	73.5	58.3	47.4
SSPSP	78.2	73.8	58.8	48.0
CPSP	78.6	74.2	59.9	48.4

can also index the i -th segment. η' is a temperature parameter controlling the concentration level of feature distribution; it is set to 0.3 in our experiments. λ'_2 is the weight to balance these two losses and is empirically set to 0.01.

The experimental result is shown in Table 9, where SSPSP is conducted with the optimal experiment setup. We can find that the performance of the SSPSP is comparable with the CPSP on the AVE dataset but is much lower on the large-scale VGGSound-AVEL100k dataset. On VGGSound-AVEL100k, our CPSP method surpasses SSPSP by 1.1% and 0.4% under fully and weakly supervised settings, respectively. This reflects that such self-supervised contrastive method is not robust for audio-visual event localization.² In fact, the self-supervised SSPSP indeed ignores the semantic alignment of audio-visual pairs when constructing positive-negative samples which is vital for AVEL. Unlike SSPSP, the proposed CPSP uses reliable audio-visual pairs to construct positive and negative samples. This makes the CPSP more superior.

6.5 Quantitative analysis - generalization to AVVP task

In this subsection, we extend the proposed CPSP method to a related and more challenging audio-visual video parsing (AVVP) task. We adopt the baseline method HAN [20] specifically designed for this task as the backbone and replace its core hybrid attention network for aggregating audio-visual features by the proposed PSP module. As for the objective optimization, we keep the loss items proposed in [20] and introduce the proposed video-level contrastive objective $\mathcal{L}_{\text{vpasa}}$ under the weakly-supervised labels (given only video-level labels) to adapt our CPSP model for AVVP.

Notably, since there are multiple categories of events in each video in AVVP task, there are some differences from AVEL to AVVP when constructing the positive and negative sets for contrastive learning. Specifically, for a certain video in a batch during training, videos in the negative set can be selected from those remaining videos that have completely

² We provide more experimental results and analyses in the appendix B.2 that show the SSPSP is much more sensitive to the data distribution and training batch size, *etc.*

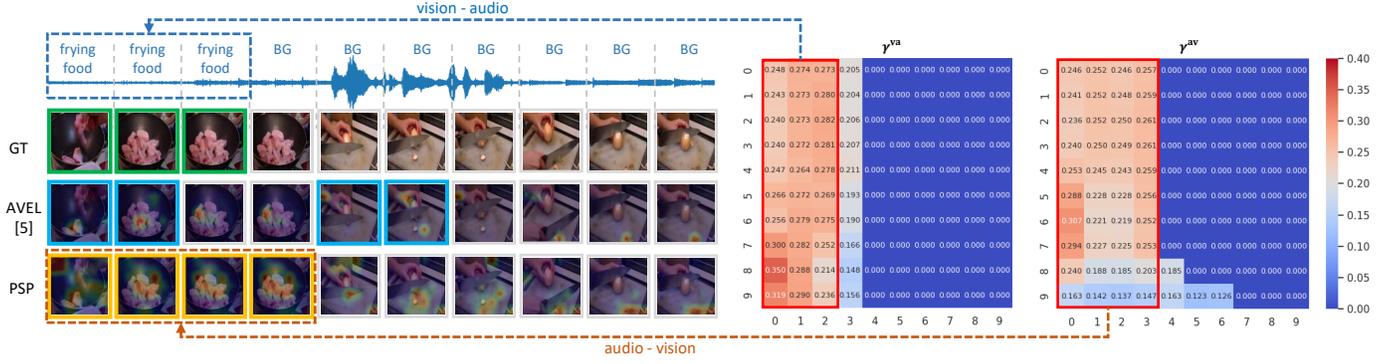


Fig. 8. A qualitative example of pair-level propagation in PSP. For the video on the **left**, only the first three segments simultaneously contain the visual and audio signals of the event *frying food*. The green boxes represent ground truth labels. The blue and orange boxes indicate predictions of AVEL [5] and the PSP method, respectively. Besides, we visualize the attention effect on the images. It is clear that our method produces more accurate localization. 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 box in γ^{va} shows that all the visual components are highly correlated with the first three audio components containing the sound of the event. 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.

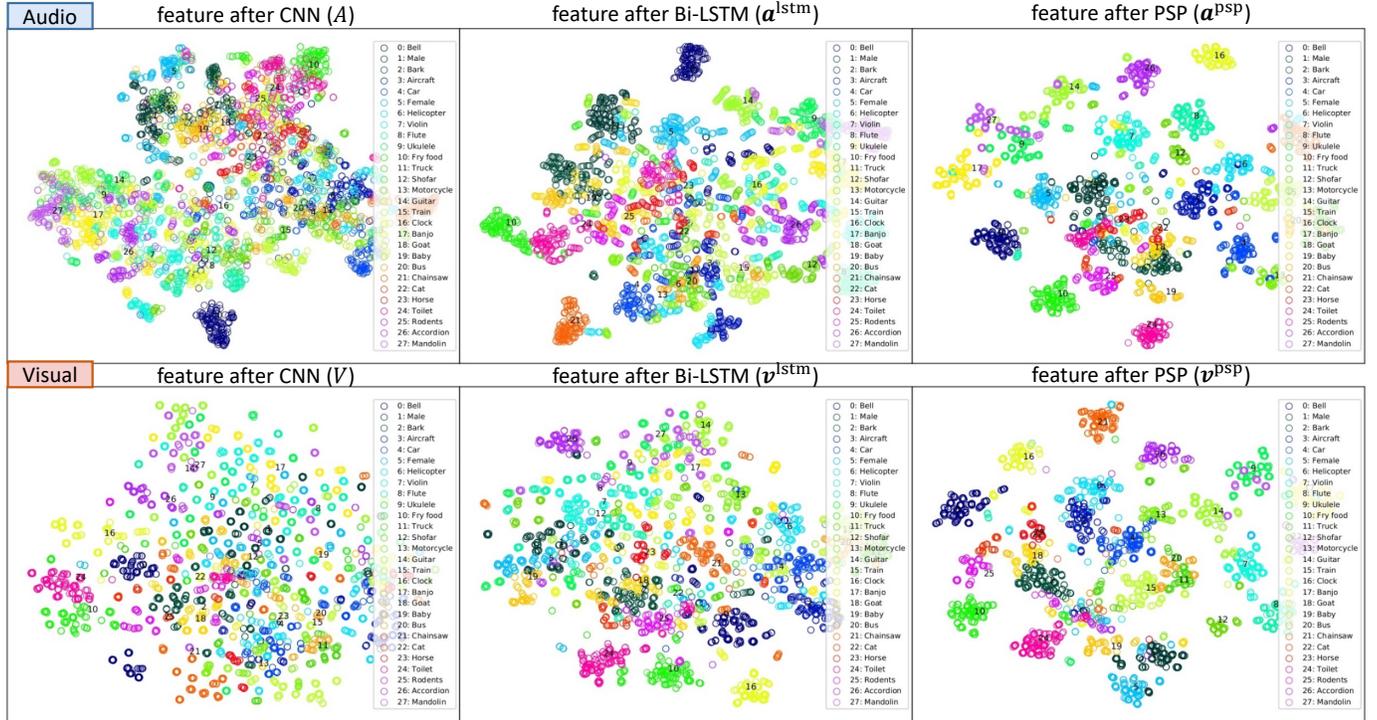


Fig. 9. TSNE [80] visualization of audio and visual feature distributions. The data all come from the validation set of AVE dataset under the fully supervised setting. (**Row 1:**) audio features. (**Row 2:**) visual features. (**Column 1:**) the CNN features. (**Column 2:**) features after Bi-LSTM encoding. (**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.

irrelevant event categories from it. But it is hard to select the positive samples requiring completely the same labels. Therefore, we consider to define a co-occurrence ratio r to measure the coincidence degree of event categories between two videos. For example, given a video vid_a with event labels $\{barking, speech\}$ and its pair video vid_b with labels $\{speech, music, clapping\}$, r is calculated by the proportion of co-occurrence event labels ($\{speech\}$) to the total event labels of vid_a ($\{barking, speech\}$), i.e., r is 1/2. In this way, the ratio r indicates that the positive samples are expected to contain as many events as possible that are appeared in the reference video (large r). We consider to set a threshold μ to construct the positive samples. For any pair of videos, we

first compute the ratio r between them. If the r is greater or equal than pre-set μ ($r \geq \mu$), the pairwise video is selected as positive sample. As for the negative samples, the ratio r is strictly equal to zero which means there are no events overlapping between the two videos. The hyper-parameters μ and θ , K in Eq. 6 are empirically set to 0.6, 0.4 and 4 for AVVP, respectively. And we provide an ablation study on the threshold μ in the appendix C.

With the above setup, we train our CPSP model on LLP dataset [20] from scratch for AVVP. For fair comparison, we use the same evaluation metrics as in HAN [20], referring to “A” and “V” (the F-score of audio events and visual events, respectively), “AV” (the F-score of audio-visual co-

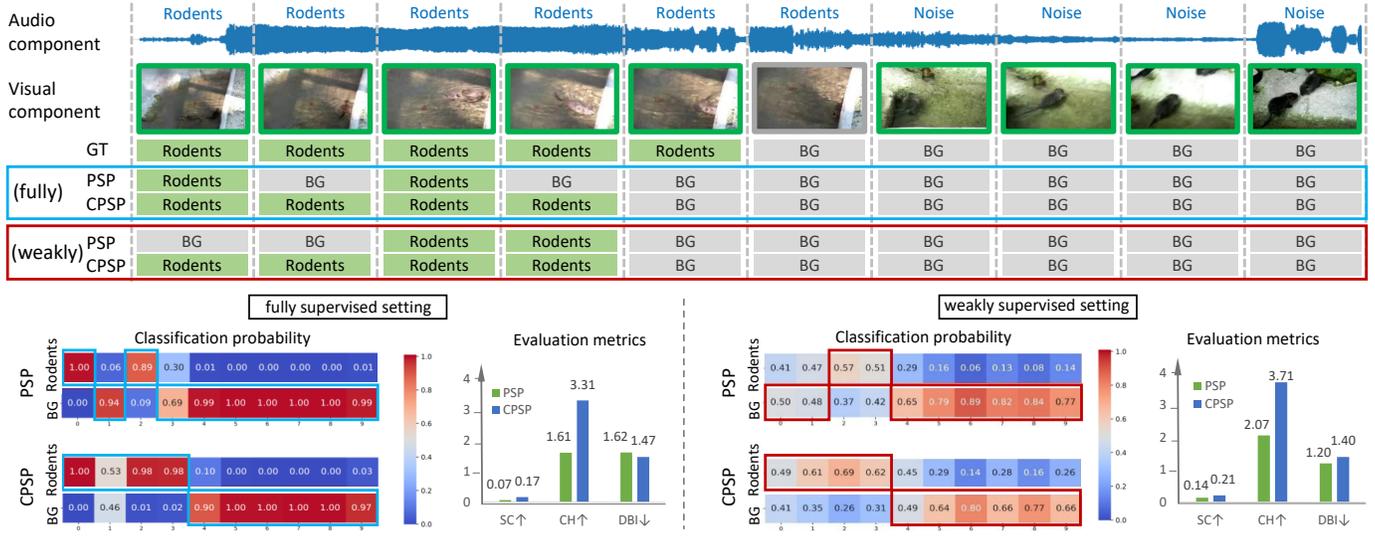


Fig. 10. Localization results of a qualitative example (sampling from \mathcal{D}_{bg} of AVE dataset) under both fully and weakly supervised settings. In this example, the first five video segments contain the audio-visual event *Rodents*. In either setting, PSP wrongly predicts some segments as *Background* (gray boxes), the CPSP method has the correct results (green boxes). The classification probability maps confirm this. We also compare the features learned by the PSP and CPSP using the metrics aforementioned in Table 7. These metrics reflect the clustering quality of features. The results show that the features learned by the CPSP are more distinguishable. It also verifies the segment-level event-awareness capability of the CPSP.

occurrence events, namely AVEs in AVEL task), “Type@AV” (the averaged result of “A”, “V”, and “AV”), and “Event@AV” (the F-score of audio-visual events where mIoU is set to 0.5). From Table 8, we have three observations: **first**, our PSP and CPSP methods surpass the baseline methods from audio-visual event localization task (AVEL [5], AVSDN [6]) by a large margin for audio-visual video parsing (AVVP). **Second**, compared to the vanilla PSP, the proposed CPSP with the video-level contrastive objective \mathcal{L}_{vpsa} improves the performances significantly. For example, the metric “Type@AV” and “Event@AV” are improved by 3.9% and 3.3% for the segment-level, and they are 4.1% and 4.3% for the event-level, respectively. This again demonstrates the benefits of the proposed contrastive strategy. **Third**, compared to the HAN [20] that is specially designed for AVVP, CPSP is even more superior that has better performances, especially having obvious performance superiority on “V” and “AV”. This reflects that the CPSP not only keeps a strong ability to recognize the audio-visual events but can also competently identify separate audio events and visual events.

In short, the proposed CPSP is still effective and advanced in both the AVEL and more challenging AVVP tasks providing more evidence of the model generalization.

6.6 Qualitative analysis

6.6.1 Visualization of the effectiveness of PSP

Propagated audio and visual components in PSP. We start by presenting an example of audio-visual event localization in Fig. 8. The event in this sample is difficult to predict because the visual images are changeable and the audio signals are mixed with background noise. From the figure, we have three observations. (1). While both our method and AVEL [5] use the AVGA attention, we show that the PSP enables better attention to visual regions closely related to sound sources. As displayed in Fig. 8, 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 decisions 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. 8. 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, they tend to build strong connections (large similarity values) between the first three audio components and the first four visual components containing the scene of “flying food”, *i.e.*, feature propagation merely takes place in these event-related audio-visual pairs. 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 back-propagation, it allows the model to be able to attend on more sound-relevant visual regions and visual-relevant audio segments.

Feature distribution in the PSP. We then visualize the data distribution of features processed by different stages in our framework using TSNE [80]. As shown in Fig. 9, 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 *Frying 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



Fig. 11. Weakly supervised setting is a more challenging setting. We display the AVE localization results of two examples that all of the video segments contain the event (sampling from \mathcal{D}_{ave} of AVE dataset). We have two observations: 1) as shown in example (a), the PSP classifies incorrect event category in the orange box, while the CPSP provides accurate predictions in green boxes; 2) for example (b), even both the PSP and CPSP have exact predictions, the CPSP gives larger probabilities to the ground truth category. The classification probability maps confirm this. This indicates that the features encoded by the CPSP contain more category-aware semantics related to the ground truth thus facilitates the event category classification.

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.

6.6.2 Visualization of the effectiveness of CPSP

Here, we display some examples to explore the classification capability of the CPSP, where compared with the PSP, the CPSP introduces the contrastive constraints PSA_S and PSA_V .

Segment-level event-aware. First, in Fig. 10, we show a video example. We perform the PSP and CPSP in both fully and weakly supervised settings and report the AVE localization results. The CPSP has more accurate predictions in both settings. Using the PSP, some segments containing the event are incorrectly classified to the background (*i.e.*, the second and fourth segments in the fully supervised setting, the first two segments in the weakly supervised setting) while the CPSP outputs all the correct results. We display the classification probability map to see what happens. In the fully supervised setting, the second and fourth segments are predicted to the *background* with high probabilities by the PSP but this result is overturned by the CPSP. Similar phenomenon is observed in the weakly supervised setting, which performs slightly lower probabilities than full supervision due to its poor knowledge (segment-level event labels). We also use the evaluation metrics mentioned in Table 7 to test the discriminability of the segment features. As shown in the histograms, the CPSP has superior performances under all of the indicators in both settings. This reflects that positive event-aware semantics of segment features are aggregated and discriminative from the backgrounds. We speculate this

is attributed to the PSA_S that enforces the CPSP to learn more *event-aware* semantics.

Video-level category-aware. We further display two examples containing no backgrounds and conduct the localization under the more challenging weakly supervised setting. As shown in Fig. 11, for example (a), the PSP incorrectly classify the seventh segment to the *Helicopter* (orange box; the ground truth is *Aircraft*); it's even easily confused by human. CPSP takes efforts to modify the wrongly predicted result generated by the PSP (orange bounding box). This is reflected in the classification probability map with a high probability of *Aircraft* at the seventh segment. In example (b), both the PSP and CPSP provide accurate predictions, *i.e.*, all the segments are classified to the *Car* event. But the classification probability map tells that the CPSP gives higher probabilities to the *Car* category (red bounding box), making the video segments more recognizable from the similar *Truck* or *background*. These two examples demonstrate the CPSP is more *category-aware* thanks to the PSA_V that enables to encode features including more semantics related to the ground truth category thus distinguishing from other categories.

7 CONCLUSION

For the AVE localization problem, we propose a contrastive positive sample propagation (CPSP) method that comprehensively explores three levels of positive samples for distinguishable audio-visual representation learning. Specifically, the pair-level PSP identifies and exploits the most relevant audio and visual samples when fusing the cross-modal features. We find that negative and weak connections, even though

with small weights, have a detrimental effect on the system, and thus need to be completely removed. The segment-level PSA_S and video-level PSA_V provide additional contrastive constraints to refine the features encoded by the PSP. The PSA_S enforces the model to be event-aware by gathering the positive segments containing an AVE and being far away from the backgrounds. The PSA_V is actually an online hard sample learning that contrasts the positive video from negatives according to the event category thus makes the model to be category-aware. We show that such pair-level, segment-level and video-level positive sample propagation and activation method are beneficial to the classifier training. To evaluate the model generalization ability, we collect a large-scale VGGSound-AVEL100k dataset and extend our method to a more challenging audio-visual video parsing task. Extensive experimental results validate the effectiveness of the proposed CPSP method. In addition, this paper covers a comprehensive study on the contrastive learning manners with different supervisions, *i.e.*, fully-, weakly-, and self-supervised.

ACKNOWLEDGMENTS

We would like to thank the reviewers for their constructive suggestions. This work was supported by the National Natural Science Foundation of China (72188101, 61725203, 62020106007, and 62272144), and the Major Project of Anhui Province (202203a05020011).

REFERENCES

- [1] R. Arandjelovic and A. Zisserman, "Look, listen and learn," in *ICCV*, 2017, pp. 609–617.
- [2] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *arXiv preprint arXiv:2103.00020*, 2021, pp. 1–36.
- [3] X. Cheng, Y. Zhong, Y. Dai, P. Ji, and H. Li, "Noise-aware unsupervised deep lidar-stereo fusion," in *CVPR*, 2019, pp. 6339–6348.
- [4] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig, "Scaling up visual and vision-language representation learning with noisy text supervision," in *ICML*, 2021, pp. 1–14.
- [5] Y. Tian, J. Shi, B. Li, Z. Duan, and C. Xu, "Audio-visual event localization in unconstrained videos," in *ECCV*, 2018, pp. 247–263.
- [6] Y.-B. Lin, Y.-J. Li, and Y.-C. F. Wang, "Dual-modality seq2seq network for audio-visual event localization," in *ICASSP*, 2019, pp. 2002–2006.
- [7] Y. Wu, L. Zhu, Y. Yan, and Y. Yang, "Dual attention matching for audio-visual event localization," in *ICCV*, 2019, pp. 6292–6300.
- [8] H. Xu, R. Zeng, Q. Wu, M. Tan, and C. Gan, "Cross-modal relation-aware networks for audio-visual event localization," in *ACM MM*, 2020, pp. 3893–3901.
- [9] H. Xuan, Z. Zhang, S. Chen, J. Yang, and Y. Yan, "Cross-modal attention network for temporal inconsistent audio-visual event localization," in *AAAI*, 2020, pp. 279–286.
- [10] J. Ramaswamy, "What makes the sound?: A dual-modality interacting network for audio-visual event localization," in *ICASSP*, 2020, pp. 4372–4376.
- [11] J. Ramaswamy and S. Das, "See the sound, hear the pixels," in *WACV*, 2020, pp. 2970–2979.
- [12] J. S. Chung and A. Zisserman, "Out of time: automated lip sync in the wild," in *ACCV*, 2016, pp. 251–263.
- [13] B. Korbar, D. Tran, and L. Torresani, "Cooperative learning of audio and video models from self-supervised synchronization," in *NeurIPS*, 2018, pp. 7774–7785.
- [14] N. Khosravan, S. Ardeshtir, and R. Puri, "On attention modules for audio-visual synchronization," in *CVPR Workshops*, 2019, pp. 1–4.
- [15] T. Afouras, A. Owens, J. S. Chung, and A. Zisserman, "Self-supervised learning of audio-visual objects from video," in *ECCV*, 2020, pp. 1–23.
- [16] S. Ma, Z. Zeng, D. McDuff, and Y. Song, "Contrastive learning of global and local audio-visual representations," in *NeurIPS*, 2021, pp. 1–11.
- [17] Y. Wu and Y. Yang, "Exploring heterogeneous clues for weakly-supervised audio-visual video parsing," in *CVPR*, 2021, pp. 1326–1335.
- [18] J. Zhou, L. Zheng, Y. Zhong, S. Hao, and M. Wang, "Positive sample propagation along the audio-visual event line," in *CVPR*, 2021, pp. 8436–8444.
- [19] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, "Vggsound: A large-scale audio-visual dataset," in *ICASSP*, 2020, pp. 721–725.
- [20] Y. Tian, D. Li, and C. Xu, "Unified multisensory perception: Weakly-supervised audio-visual video parsing," in *ECCV*, 2020, pp. 436–454.
- [21] R. Arandjelovic and A. Zisserman, "Objects that sound," in *ECCV*, 2018, pp. 435–451.
- [22] Y. Aytar, C. Vondrick, and A. Torralba, "Soundnet: Learning sound representations from unlabeled video," in *NeurIPS*, 2016, pp. 1–9.
- [23] Y. Cheng, R. Wang, Z. Pan, R. Feng, and Y. Zhang, "Look, listen, and attend: Co-attention network for self-supervised audio-visual representation learning," in *ACM MM*, 2020, pp. 3884–3892.
- [24] H. M. Fayek and A. Kumar, "Large scale audiovisual learning of sounds with weakly labeled data," in *IJCAI*, 2020, pp. 558–565.
- [25] T. Darrell, J. W. Fisher, and P. Viola, "Audio-visual segmentation and "the cocktail party effect"," in *International Conference on Multimodal Interfaces*, 2000, pp. 32–40.
- [26] C. Xu, W. Rao, X. Xiao, E. S. Chng, and H. Li, "Single channel speech separation with constrained utterance level permutation invariant training using grid lstm," in *ICASSP*, 2018, pp. 6–10.
- [27] T. Afouras, J. S. Chung, A. Senior, O. Vinyals, and A. Zisserman, "Deep audio-visual speech recognition," in *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2018, pp. 1–13.
- [28] T. Jenrungrot, V. Jayaram, S. Seitz, and I. Kemelmacher-Shlizerman, "The cone of silence: speech separation by localization," in *NeurIPS*, 2020, pp. 1–17.
- [29] R. Gao and K. Grauman, "Visualvoice: Audio-visual speech separation with cross-modal consistency," in *CVPR*, 2021, pp. 15 495–15 505.
- [30] S. Parekh, S. Essid, A. Ozerov, N. Q. Duong, P. Pérez, and G. Richard, "Motion informed audio source separation," in *ICASSP*, 2017, pp. 6–10.
- [31] J. Pu, Y. Panagakis, S. Petridis, and M. Pantic, "Audio-visual object localization and separation using low-rank and sparsity," in *ICASSP*, 2017, pp. 2901–2905.
- [32] H. Zhao, C. Gan, A. Rouditchenko, C. Vondrick, J. McDermott, and A. Torralba, "The sound of pixels," in *ECCV*, 2018, pp. 570–586.
- [33] R. Gao, R. Feris, and K. Grauman, "Learning to separate object sounds by watching unlabeled video," in *ECCV*, 2018, pp. 35–53.
- [34] H. Zhao, C. Gan, W.-C. Ma, and A. Torralba, "The sound of motions," in *ICCV*, 2019, pp. 1735–1744.
- [35] R. Gao and K. Grauman, "Co-separating sounds of visual objects," in *ICCV*, 2019, pp. 3879–3888.
- [36] R. Qian, D. Hu, H. Dinkel, M. Wu, N. Xu, and W. Lin, "Multiple sound sources localization from coarse to fine," in *ECCV*, 2020, pp. 1–16.
- [37] D. Hu, F. Nie, and X. Li, "Deep multimodal clustering for unsupervised audiovisual learning," in *CVPR*, 2019, pp. 9248–9257.
- [38] D. Hu, R. Qian, M. Jiang, X. Tan, S. Wen, E. Ding, W. Lin, and D. Dou, "Discriminative sounding objects localization via self-supervised audiovisual matching," *NeurIPS*, pp. 1–14, 2020.
- [39] J. Zhou, J. Wang, J. Zhang, W. Sun, J. Zhang, S. Birchfield, D. Guo, L. Kong, M. Wang, and Y. Zhong, "Audio-visual segmentation," in *ECCV*, 2022, pp. 386–403.
- [40] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *CVPR*, 2016, pp. 2921–2929.
- [41] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *ICCV*, 2017, pp. 618–626.
- [42] A. Senocak, T.-H. Oh, J. Kim, M.-H. Yang, and I. S. Kweon, "Learning to localize sound sources in visual scenes: Analysis and applications," in *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2021, pp. 1605–1619.

- [43] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, "Cnn architectures for large-scale audio classification," in *ICASSP*, 2017, pp. 131–135.
- [44] Q. Kong, Y. Xu, W. Wang, and M. D. Plumbley, "Audio set classification with attention model: A probabilistic perspective," in *ICASSP*, 2018, pp. 316–320.
- [45] A. Kumar, M. Khadkevich, and C. Fügen, "Knowledge transfer from weakly labeled audio using convolutional neural network for sound events and scenes," in *ICASSP*, 2018, pp. 326–330.
- [46] B. McFee, J. Salamon, and J. P. Bello, "Adaptive pooling operators for weakly labeled sound event detection," in *IEEE Trans. on Audio, Speech, and Language Processing*, 2018, pp. 2180–2193.
- [47] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *CVPR*, 2014, pp. 1725–1732.
- [48] X. Long, C. Gan, G. De Melo, J. Wu, X. Liu, and S. Wen, "Attention clusters: Purely attention based local feature integration for video classification," in *CVPR*, 2018, pp. 7834–7843.
- [49] X. Long, C. Gan, G. Melo, X. Liu, Y. Li, F. Li, and S. Wen, "Multimodal keyless attention fusion for video classification," in *AAAI*, 2018, pp. 1–8.
- [50] L. Wang, W. Li, W. Li, and L. Van Gool, "Appearance-and-relation networks for video classification," in *CVPR*, 2018, pp. 1430–1439.
- [51] D. Tran, H. Wang, L. Torresani, and M. Feiszli, "Video classification with channel-separated convolutional networks," in *ICCV*, 2019, pp. 5552–5561.
- [52] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," in *IEEE Trans. on Signal Processing*, 1997, pp. 2673–2681.
- [53] B. Duan, H. Tang, W. Wang, Z. Zong, G. Yang, and Y. Yan, "Audio-visual event localization via recursive fusion by joint co-attention," in *WACV*, 2021, pp. 4013–4022.
- [54] Y. Wu and Y. Yang, "Exploring heterogeneous clues for weakly-supervised audio-visual video parsing," in *CVPR*, 2021.
- [55] Y.-B. Lin, H.-Y. Tseng, H.-Y. Lee, Y.-Y. Lin, and M.-H. Yang, "Exploring cross-video and cross-modality signals for weakly-supervised audio-visual video parsing," *NeurIPS*, pp. 1–13, 2021.
- [56] Y.-B. Lin and Y.-C. F. Wang, "Audiovisual transformer with instance attention for audio-visual event localization," in *ACCV*, 2020, pp. 1–17.
- [57] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," in *arXiv preprint arXiv:1807.03748*, 2018, pp. 1–13.
- [58] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020, pp. 9729–9738.
- [59] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, 2020, pp. 1597–1607.
- [60] P. Morgado, N. Vasconcelos, and I. Misra, "Audio-visual instance discrimination with cross-modal agreement," in *CVPR*, 2021, pp. 12 475–12 486.
- [61] D. Harwath, A. Recasens, D. Surís, G. Chuang, A. Torralba, and J. Glass, "Jointly discovering visual objects and spoken words from raw sensory input," in *ECCV*, 2018, pp. 649–665.
- [62] Z. Zhang, Z. Zhao, Z. Lin, X. He *et al.*, "Counterfactual contrastive learning for weakly-supervised vision-language grounding," in *NeurIPS*, 2020, pp. 18 123–18 134.
- [63] C. Zhang, M. Cao, D. Yang, J. Chen, and Y. Zou, "Cola: Weakly-supervised temporal action localization with snippet contrastive learning," in *CVPR*, 2021, pp. 16 010–16 019.
- [64] T. Gupta, A. Vahdat, G. Chechik, X. Yang, J. Kautz, and D. Hoiem, "Contrastive learning for weakly supervised phrase grounding," in *ECCV*, 2020, pp. 752–768.
- [65] M. Ki, Y. Uh, W. Lee, and H. Byun, "In-sample contrastive learning and consistent attention for weakly supervised object localization," in *ACCV*, 2020, pp. 1–16.
- [66] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *NeurIPS*, 2020, pp. 18 661–18 673.
- [67] Z. Zeng, K. He, Y. Yan, Z. Liu, Y. Wu, H. Xu, H. Jiang, and W. Xu, "Modeling discriminative representations for out-of-domain detection with supervised contrastive learning," in *ACL*, 2021, pp. 1–9.
- [68] H. Sedghamiz, S. Raval, E. Santus, T. Alhanai, and M. Ghassemi, "Supcl-seq: Supervised contrastive learning for downstream optimized sequence representations," in *EMNLP*, 2021, pp. 1–6.
- [69] B. Gunel, J. Du, A. Conneau, and V. Stoyanov, "Supervised contrastive learning for pre-trained language model fine-tuning," in *ICLR*, 2020, pp. 1–15.
- [70] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015, pp. 1–14.
- [71] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *NeurIPS*, pp. 1097–1105, 2012.
- [72] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *ICASSP*, 2017, pp. 776–780.
- [73] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2014, pp. 1–15.
- [74] J. Yu, Y. Cheng, and R. Feng, "Mpn: Multimodal parallel network for audio-visual event localization," in *ICME*, 2021, pp. 1–6.
- [75] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 1–11.
- [76] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *AISTATS*, 2010, pp. 249–256.
- [77] P. J. Rousseeuw, "Silhouettes: A graphical aid to the interpretation and validation of cluster analysis," in *Journal of Computational and Applied Mathematics*, 1987, pp. 53–65.
- [78] T. Caliński and J. Harabasz, "A dendrite method for cluster analysis," in *Communications in Statistics-Simulation and Computation*, 1974, pp. 1–27.
- [79] D. L. Davies and D. W. Bouldin, "A cluster separation measure," in *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 1979, pp. 224–227.
- [80] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," in *JMLR*, 2008, pp. 2579–2605.

APPENDIX A

INTRODUCTION OF THE VGGSound-AVEL100K DATASET.

In this work, the newly collected dataset provides an additional option for evaluating the generalization ability of the designed AVEL models. We wish it will facilitate relevant research in audio-visual community. VGGSound-AVEL100k contains 101,072 videos covering 141 event categories, where the video are sampled from the VGGSound [19] dataset. Category names and video number for each category are provided on our project homepage <https://github.com/jasongief/CPSP>.

For data labeling, the category label can be easily obtained according to the video tags, while the temporal event labels (audio-visual correspondence labels) are manually labeled by outsourcing. Similar to the annotation process of the AVE [5] dataset, for each video in VGGSound-AVEL100k, each annotator is required to watch through the entire video and label each video segment: if one event simultaneously occurs in both audio and visual channels, the label for the current segment is "1", otherwise "0". To ensure the quality of labeling, 60% of the annotated data are randomly selected and manually checked again. In this way, we can obtain the final labels of VGGSound-AVEL100k dataset. More details such as annotation guideline, annotator information, and quality control can be seen from our website <https://github.com/jasongief/CPSP>.

APPENDIX B

EXPERIMENTS FOR AUDIO-VISUAL EVENT LOCALIZATION.

B.1 Multi-run experiments of CPSP

To evaluate the robustness of our method, we perform the proposed CPSP method on the AVE [5] and the new large-scale VGGSound-AVEL100k datasets by running five times under both the fully and weakly supervised settings. For comparison, we also reproduce some popular baseline methods - AVEL [5], AVSDN [6], CMRA [8], CMAN [9], and PSP [18]. The results are shown in Table 10 from which we have these observations: **first**, the top performances on AVE and VGGSound-AVEL100k during five runs that demonstrate the generalization ability of CPSP to both small and large-scale datasets. **Second**, in all the cases, the vanilla PSP still outperforms most of the SOTAs [5], [6], [9]. CMRA [8] is competitive with the PSP, but the proposed CPSP method beats both of them. Take the VGGSound-AVEL100k for an example, compared to the PSP, the average performance of CPSP exceeds it by 1.9% (from 57.8% to 59.7%) and 1.1% (from 47.1% to 48.2%) for the fully and weakly supervised settings, respectively. **Third**, the proposed CPSP keeps high performance but with relatively low accuracy variance (Accuracy@std) that show its robustness. Notably, to keep consist with public literature in the main paper, we report the best performance on AVE dataset and the same to VGGSound-AVEL100k dataset.

TABLE 10

Comparison of multiple runs on the AVE and VGGSound-AVEL100k datasets in both fully and weakly supervised settings. Accuracy@avg and Accuracy@std denote the averaged and the standard deviation of the accuracy values, respectively.

Dataset	Setting	Method	Run-1	Run-2	Run-3	Run-4	Run-5	Accuracy@avg \uparrow	Accuracy@std \downarrow
AVE [5]	fully	AVEL [5]	70.5	69.9	70.5	69.0	71.3	70.2	0.8
		AVSDN [6]*	70.8	72.3	72.0	71.2	69.7	71.2	0.9
		CMAN [9]*	70.7	72.3	71.7	71.2	71.3	71.4	0.5
		CMRA [8]	74.8	74.1	74.7	74.6	74.3	74.5	0.3
		PSP [18]	77.8	77.4	76.3	77.3	76.8	77.1	0.5
		CPSP	78.6	77.9	78.0	78.3	78.3	78.2	0.2
	weakly	AVEL [5]	66.1	66.2	66.0	66.3	65.0	65.9	0.5
		AVSDN [6]*	65.3	62.6	61.8	62.0	58.9	62.1	2.0
		CMAN [9]*	68.6	68.8	68.5	68.6	68.8	68.7	0.1
		CMRA [8]	71.2	71.3	70.7	71.6	71.8	71.3	0.4
		PSP [18]	73.5	72.2	73.3	72.8	72.5	72.9	0.5
		CPSP	74.2	73.6	74.0	73.8	73.9	73.9	0.2
VGGSound-AVEL100k	fully	AVEL [5]	55.7	55.1	54.8	55.6	55.6	55.4	0.3
		CMRA [8]*	57.1	56.6	56.3	56.7	56.9	56.7	0.3
		PSP [18]	58.3	57.1	57.6	57.9	58.0	57.8	0.4
		CPSP	59.9	59.7	59.8	59.5	59.5	59.7	0.2
	weakly	AVEL [5]	45.5	46.0	45.6	45.7	46.2	45.8	0.3
		CMRA [8]*	46.8	46.6	46.3	46.0	46.8	46.5	0.3
		PSP [18]	47.4	46.7	47.4	46.5	47.4	47.1	0.4
		CPSP	48.3	48.2	47.8	48.3	48.4	48.2	0.2

* indicates that method is reproduced by us. \square marks the best performance in the multi-runs.

B.2 More discussion on the comparison to self-supervised methods

Comparison to the SSPSP. As introduced in the main paper, SSPSP directly takes the synchronized audio-visual segment pair as a positive sample and otherwise is negative. For the self-supervised SSPSP, we observe the new results from three aspect as follows.

1) *SSPSP uses a huge order of magnitude more <positive, negative> samples* than CPSP for contrastive learning. Given a batch of video data, we denote the number of videos in this batch as N^b . In SSPSP, for each audio segment, all the visual segments from the batch except the synchronized segment (*including all the intra- and inter- video segments in a batch*) are treated as negative samples. The number of negative samples is $N^b T \times (N^b T - 1)$, where T denotes each video contains T video segments. In contrast, the CPSP merely adopts *background segments in a video* as negative samples. Therefore, the negatives used in CPSP is just $N^b \times (T - n^{bg}) \times n^{bg}$, where n^{bg} is the number of background segments per video. There is $(N^b T - 1) \gg n^{bg}$. Obviously, the amount of negative samples used in the SSPSP is a huge order of magnitude more than that in CPSP. For example, by statistics, when extending to the whole dataset, the number of negative samples at segment-level used during training of SSPSP and CPSP is approximately 1537 : 1 on the AVE dataset. As a result, SSPSP achieves promising performance that attributes the success to data-driven with numerous negative samples. However, this cannot completely dispel the doubt that SSPSP sounds not solidly and technically for the AVE Localization task. The doubt is that SSPSP will inevitably bring false negatives of audio-visual pair depicting the same event but at different timestamp. Theoretically, these audio-visual pairs are audio-visual correspondence but will be wrongly considered as the negatives.

2) *SSPSP benefits from subset \mathcal{D}_{bg} but distorts from \mathcal{D}_{ae} in each dataset.* We examine the methods on the data distribution. As shown in Table 11, SSPSP performs much better on the video type \mathcal{D}_{bg} than \mathcal{D}_{ae} . \mathcal{D}_{bg} denotes the video that contains both AVE and background segments (as shown in Fig. 5(a)); \mathcal{D}_{ae} represents another type of video that contains pure AVE segments belonging to a certain event category (as shown in Fig. 5(b)). For example, on the large-scale VGGSound-AVEL100k dataset, the SSPSP has terrible performance drop compared with the proposed CPSP from 58.8% \rightarrow 50.2% \downarrow with \mathcal{D}_{ae} . It reflects the weakness of such self-supervised learning in theory - taking the unsynchronized audio-visual segment pairs but still semantic-corresponding (positive samples in fact) as the negatives introduces extra noises for the AVE localization task. The performance of SSPSP obviously drops in the case of fine-tuning with \mathcal{D}_{ae} compared to \mathcal{D}_{bg} . Unlike this, the CPSP performs stably and is almost not affected under all the conditions.

3) *SSPSP is sensitive to the training batch size.* As well known, the performance of such self-supervised method is influenced by the number of the negative samples during training. We test the influence of batch size here. As shown in Table 12, the performance of the SSPSP obviously decreases under the small batch size (58.8% with batch size 128 \rightarrow 55.8% with batch size 32). Unlike this, the CPSP keeps highly stable.

To summarize, such self-supervised learning can bring some improvements for AVEL by learning from enormous audio-visual pairs but it is much more sensitive to the data distribution and the training batch size. Facing these factors, the CPSP is still superior and robust.

TABLE 11

Evaluation on each subset for the CPSP and the SSPSP in the fully supervised setting, measured by accuracy(%) on the AVE and VGGSound-AVEL100k datasets.

Method	AVE ($\#\mathcal{D}_{ae} : \#\mathcal{D}_{bg} \approx 2:1$)			VGGSound-AVEL100k ($\#\mathcal{D}_{ae} : \#\mathcal{D}_{bg} \approx 3:2$)		
	\mathcal{D}_{ae}	\mathcal{D}_{bg}	$\mathcal{D}_{ae} \& \mathcal{D}_{bg}$	\mathcal{D}_{ae}	\mathcal{D}_{bg}	$\mathcal{D}_{ae} \& \mathcal{D}_{bg}$
PSP	-	-	77.8	-	-	58.3
SSPSP	77.4	77.9	78.2	50.2	57.6	58.8
CPSP	78.2	78.3	78.6	58.6	59.2	59.9

* $\#\mathcal{D}_{bg}$ denotes the size of subset \mathcal{D}_{bg} , and the same is $\#\mathcal{D}_{ae}$.

TABLE 12

Performances of the CPSP and SSPSP with different training batch-sizes on the VGGSound-AVEL100k dataset in the fully supervised setting, measure by the accuracy(%).

Method	Batch size		
	128	64	32
SSPSP	58.8	57.7	55.8
CPSP	59.9	59.7	59.5

Comparison to the self-supervised method Global-Local [16]. Global-Local [16] is built upon the large-scale pre-training and achieves remarkable performance on the downstream audio-visual event localization task. It performs both spatial and temporal audio-visual contrastive learning in a self-supervised manner. To compare with Global-Local [16], we adopt the same visual backbone, *i.e.*, the advanced 3D-ResNet, to extract the visual features. The results are shown in Table 13, our method is comparable to the Global-Local [16]. For the fully supervised setting, the proposed CPSP with segment-level and video-level positive sample activation (PSA_S and PSA_V) achieves better performance than Global-Local [16]

(82.3% vs. 82.1%). As for the weakly supervised setting, the CPSP equipped with only video-level PSA_V is comparable to Global-Local [16] (78.7% vs. 79.8%).

But importantly, our method achieves such comparable performances at a much lower cost. 1) The Global-Local [16] is pretrained on an extra large-scale K-AV-240K dataset that contains 240k audio-visual videos, while our method does not need any extra data. 2) To achieve better performances, the Global-Local [16] samples video frames at 10 FPS on AVE [5] dataset for data augmentation. In our work, in order to keep the consistent experiment setup, we keep 1 FPS in our method as the same as existing literature [5]–[11], [74]. 3) The authors of Global-Local [16] suggest using 16 Tesla P100 GPU to handle the large-scale dataset of 240k videos while our model can be trained very lightly with just one GTX 1080 GPU without extra data. Consequentially, it is obvious that our method spends less training time and fewer GPUs since it uses much fewer data.

To summarize, even without those techniques used in the Global-Local [16] (*e.g.*, large-scale pre-training, dense video frames sampling, high-performance GPU device), our model still performs competitively.

TABLE 13

Comparison with the self-supervised Global-Local built upon large-scale pre-training [16] for audio-visual event localization.

Method	Dataset		#GPU	Visual Encoder	Setting	
	extra data	train data			fully	weakly
CPSP	none	AVE [5]]	1 GTX 1080	VGG-19	78.6	74.2
				3D-ResNet	82.3	78.7
Global-Local [16]	K-AV-240K [16]	AVE [5]	16 Tesla P100	3D-ResNet	82.1	79.8

APPENDIX C

EXPERIMENTS FOR AUDIO-VISUAL VIDEO PARSING.

Ablation study on the threshold μ . As introduced in the main paper, the threshold μ is used to select the positive samples for the AVVP task. We perform an ablation study on μ and the results are shown in Table 14. The best performance of CPSP is achieved when μ is set to 0.6. When μ is set to 1.0, it means only videos have the exactly same label set of multiple instance categories will be selected as positives. This condition is kind of strict especially in a batch video data during training. We argue that $\mu = 0.6$ (60% category coverage) is a suitable parameter to make the model flexible in selecting positive samples during training and we choose this setup in the experiments for AVVP. We have also released the pretrained model at <https://drive.google.com/drive/folders/1cqlVeAKx1NFKnk0ynXvyyQHnnKqyecNm?usp=sharing>.

TABLE 14

Ablation study on the threshold μ used for $\mathcal{L}_{\text{vpsa}}$ to construct positive and negative samples for AVVP task. Experiments are conducted on the LLP dataset.

μ	Segment-level					Event-level				
	A	V	AV	Type@AV	Event@AV	A	V	AV	Type@AV	Event@AV
0.5	58.2	56.3	51.4	55.3	55.5	51.0	52.4	45.3	49.6	49.3
0.6	58.5	57.8	52.6	56.3	55.8	51.6	54.0	46.5	50.7	49.9
1.0	56.2	56.8	50.2	54.4	54.8	48.9	52.3	44.0	48.4	48.0