

# Self-Supervised Predictive Learning: A Negative-Free Method for Sound Source Localization in Visual Scenes

Zengjie Song<sup>1</sup> Yuxi Wang<sup>1,3</sup> Junsong Fan<sup>1,2</sup> Tieniu Tan<sup>1,2</sup> Zhaoxiang Zhang<sup>1,2,3\*</sup>

<sup>1</sup>Center for Research on Intelligent Perception and Computing, NLPR, CASIA

<sup>2</sup>University of Chinese Academy of Sciences (UCAS)

<sup>3</sup>Centre for Artificial Intelligence and Robotics, HKISI-CAS

{zengjie.song, wangyuxi2016, fanjunsong2016, zhaoxiang.zhang}@ia.ac.cn, tnt@nlpr.ia.ac.cn

## Abstract

Sound source localization in visual scenes aims to localize objects emitting the sound in a given image. Recent works showing impressive localization performance typically rely on the contrastive learning framework. However, the random sampling of negatives, as commonly adopted in these methods, can result in misalignment between audio and visual features and thus inducing ambiguity in localization. In this paper, instead of following previous literature, we propose Self-Supervised Predictive Learning (SSPL), a negative-free method for sound localization via explicit positive mining. Specifically, we first devise a three-stream network to elegantly associate sound source with two augmented views of one corresponding video frame, leading to semantically coherent similarities between audio and visual features. Second, we introduce a novel predictive coding module for audio-visual feature alignment. Such a module assists SSPL to focus on target objects in a progressive manner and effectively lowers the positive-pair learning difficulty. Experiments show surprising results that SSPL outperforms the state-of-the-art approach on two standard sound localization benchmarks. In particular, SSPL achieves significant improvements of 8.6% cIoU and 3.4% AUC on SoundNet-Flickr compared to the previous best. Code is available at: <https://github.com/zjsong/SSPL>.

## 1. Introduction

When strolling in a park brimming with life, you notice that the bird sitting on a twig is chirping; the puppy on the road ahead gives a little bark; and after a while an acquaintance may walk by and say hello to you friendly. Despite a short notice, humans own the excellent ability to associate

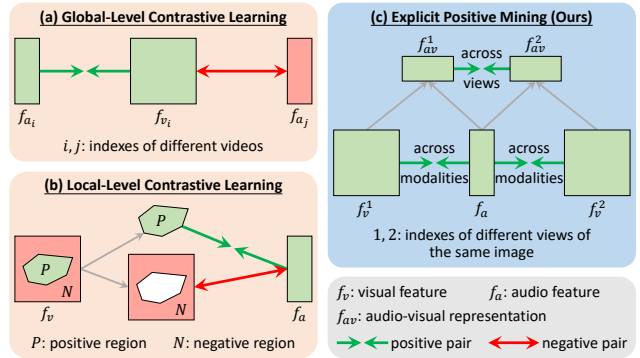


Figure 1. **Three methods to build audio-visual correspondence for sound localization.** Previous contrastive learning based methods have to construct *negative* pairs (a) at the global feature level or (b) at the local feature level. (c) Our method, by contrast, explores the coherent similarities between audio and visual features by only explicit *positive* mining.

the sounds they hear with the corresponding visual perception, and thus can localize and distinguish different sounding objects from one another.

To mimic humans' such ability, in this work, we pay attention to the task of sound source localization in visual scenes, where the goal is to localize regions of the visual landscape that correlate highly with the audio cues. While handling this task is a long-standing challenge [17, 25], remarkable breakthroughs have been made until recent progresses on self-supervised audio-visual learning [1, 27, 28, 41, 43, 46]. These methods leverage the free supervision rooted in videos, *e.g.*, the natural correspondence and/or temporal synchronization between audio and visual sources, to guide multi-modal feature extraction and alignment; then the similarity map between audio and visual features is usually employed to localize sounding objects. Among them, contrastive learning has particularly achieved impressive performance on this task [1, 5, 43, 46, 51].

Existing contrastive learning methods in this line of work

\*Corresponding author.

can be cast into two categories: the first one is global-level contrastive learning (GLCL, Figure 1a) [1, 37, 43, 46], which commonly attracts audio and visual features extracted from the same video and repulses features from different videos; the other one is local-level contrastive learning (LLCL, Figure 1b) [5, 35, 51], which further compares audio feature with different visual feature components, even though they have correspondence at the video level. Generally, to perform contrastive learning, these methods randomly sample sounds to form negative pairs with the given video frame. However, this randomness can produce false negatives by sampling sounds that actually belong to the same category as the positive one, and thus hampering the model to align audio and visual features in semantic level. The misalignment as a result induces the learning process to build inaccurate audio-visual correspondence for localization.

We conduct a pilot experiment to illustrate the effect of such false negatives in Figure 2. Given image and sound from the same video (*e.g.*, saxophone playing) as positive pair, other videos’ sounds, holding the same category as the positive one, are allowed to construct negative pairs in Figure 2a, while not allowed in Figure 2b. We keep the remaining training settings same for these two cases. During testing, consequently, the former case generates ambiguous localization on sounding objects (*i.e.*, saxophone here) and the later one not. Based on this observation, we take a step back and ask the questions: Do we really need negatives to develop self-supervised sound localization methods? Can the image-audio positive pair alone be used to achieve the same goal?

To answer these questions, we propose Self-Supervised Predictive Learning (SSPL), a negative-free approach for sound source localization through explicit audio-visual positive mining. The predictive learning is embodied from two perspectives (Figure 1c): predicting across different visual views and predicting across audio and visual modalities. Given an image-audio pair from one video, the former perspective hypothesizes that if two different visual views of one video frame contain the same sounding objects, they should share a consistent correspondence with the given sound source. We achieve this consistency by the mutual prediction of two audio-visual representations. For the later perspective, we devise a predictive coding module (PCM) that uses visual features to iteratively predict audio ones, providing a coarse-to-fine way to automatically align features. To the best of our knowledge, this is the first attempt to apply the self-supervised negative-free method to the audio-visual task of sound localization.

Our main contributions can be summarized as follows:

- We propose a novel negative-free method to extend a self-supervised learning framework to the audio-visual data domain for sound localization, and show how it can effectively address the false negative sampling

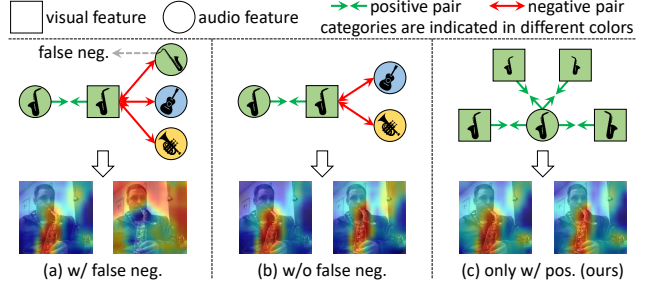


Figure 2. **Effect of false negatives on sound localization.** (a) Ambiguity in localization is observed when learning with false negatives that have the same class label as the positive one. (b) Consistent localization can be obtained without sampling false negatives, but requiring *class label* as guidance. (c) Our method mitigates this problem via self-supervised positive mining. Experiments are performed on MUSIC [56] where labels are available.

problem.

- We propose the predictive coding module for feature alignment, which enables the model to progressively attend to relevant visual features while ignoring information irrelevant to audio cues, boosting sound localization significantly.
- Comprehensive experiments demonstrate the effectiveness of the proposed approach, which achieves localization performance superior to the state-of-the-art on SoundNet-Flickr and VGG-Sound Source.

## 2. Related Work

**Self-Supervised Visual Representation Learning.** Self-supervised learning (SSL) has achieved remarkable breakthroughs on large computer vision benchmarks. Most of the current SSL methods [7, 9–11, 16, 22, 23, 44, 52] resort to the design of contrastive learning strategy [40]. These methods, at their core, transform one image into multiple views, and repulse different images (negatives) meanwhile attracting the same image’s different views (positives). Recently, several efforts have been made to further relieve the requirement of negatives and simplify the SSL framework beyond conventional contrastive learning, including BarlowTwins [55], W-MSE [15], BYOL [21], and SimSiam [12]. In SimSiam [12], researchers investigate the importance of simple Siamese architecture for unsupervised representation learning, and empirically show that the stop-gradient operation is critical for the network to prevent collapse, even without using momentum encoder [22] and large batches [9]. These advances in image representation learning provide insights for our work to develop effective audio-visual SSL method.

**Audio-Visual Representation Learning.** The vision and sound are usually two co-occurring modalities, which can naturally be used to derive supervisions for audio-visual

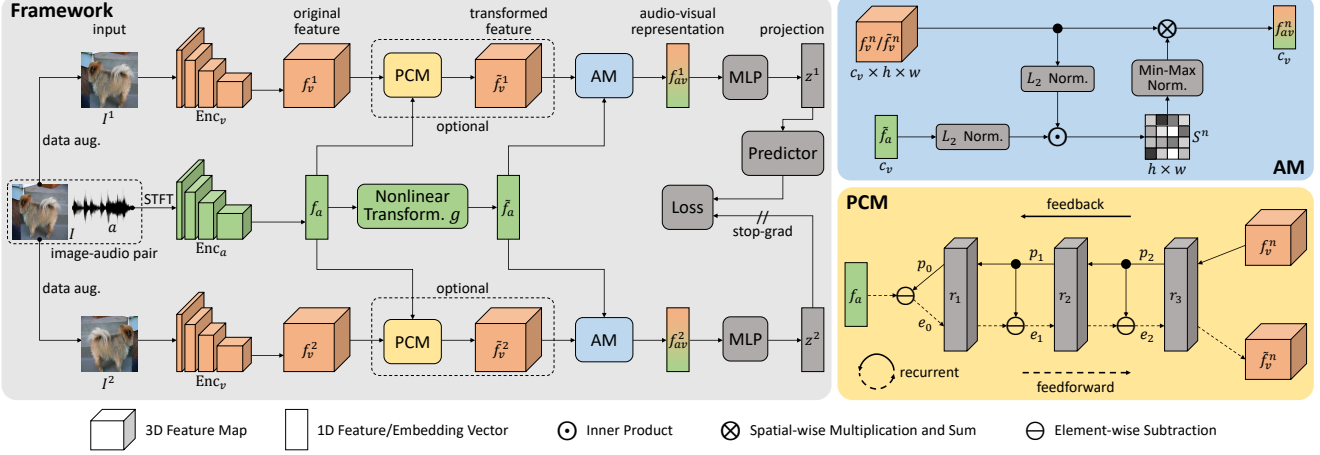


Figure 3. **Left:** framework of our SSPL method. **Top right:** attention module (AM) to compute audio-visual representation. **Bottom right:** predictive coding module (PCM) to align audio and visual features (for simplicity we only show a 3-layer version). In the framework,  $\text{Enc}_a$  of the middle processing stream derives the discriminative audio feature of audio signal  $a$ , while  $\text{Enc}_v$  extracts visual features of two augmented images  $I^1$  and  $I^2$ . Note that  $\text{Enc}_v$  shares weights between top and bottom streams, and a similar setting holds for PCM and MLP. During inference, the similarity map  $S^n$ , generated in AM, is resized to image scale and treated as the sound localization map.

learning [4, 8, 33, 42, 50]. In [4], for instance, the visual features extracted from pretrained teacher networks act to guide the student network to learn more discriminative sound representation, and vice versa in [42]. Korbarr *et al.* [33] and Owens and Efros [41] leverage the synchronization between audio and visual streams to build negative samples and contrastive losses, obtaining versatile multisensory features, respectively. Several works also explore the audio-visual correspondence by feature clustering [2, 27, 29]. In general, these methods focus on learning task-agnostic representations, which work well on classification-related down-stream tasks, such as action/scene recognition [2, 13, 33, 37, 41], audio event classification [2, 3, 27, 33, 38], video retrieval [13], *etc.* However, they are not customized for sound source localization, and as a result only achieve limited performance on this task [27, 29, 41].

**Sound Localization in Visual Scenes.** Early works to solve this task mainly rely on statistical modeling of the cross-modal relationship by using, for example, mutual information [17, 25] and canonical correlation analysis [30, 32]. However, these methods as shallow models only show advantages in simple audio-visual scenarios. By digging into the correspondence between deep audio and visual features, recent deep learning methods give promising solutions to this problem [1, 5, 18, 27, 43, 46, 51, 53]. For instance, Senocak *et al.* [46] employ a two-stream framework and an attention mechanism to compute sound localization map. Qian *et al.* [43] achieve the same goal by using the class activation map (CAM) derived from a weakly-supervised approach. In [27, 29], audio and visual features are clustered, respectively, and the assignment weights based on the distance between features and cluster centers are adopted to

localize sounding objects. In addition to viewing image and sound extracted from different videos as negative pair, Chen *et al.* [5] and Lin *et al.* [35] propose to mine hard negatives within an image-audio pair, *i.e.*, background regions that correlate lowly with the given sound are treated as extra hard negatives. Different from these negative-based works, we handle the same task by explicit positive mining, providing an effective alternative for sound localization.

### 3. Method

Figure 3 depicts the overall framework of our SSPL, which is a three-stream network, making a big difference with widely-used two-stream ones. The top and bottom streams first serve to extract deep visual features from different views of the same image. Then, they employ PCM and AM to integrate visual and audio features, where the discriminative audio feature is derived from the middle stream subnetwork. Subsequently, two audio-visual representations are enforced to be similar by self-supervised predicting with each other. The sound localization map is a natural consequence of representation learning and is generated in the AM. Note that the vanilla SSPL without PCM focuses on exploring audio-visual correspondence across different image views (Sec. 3.2), while the PCM component excels at aligning features across modalities (Sec. 3.3), and thus boosting localization performance further. We elaborate on and formulate each part in the following.

#### 3.1. Unimodal Features of Audio and Vision

Let  $I \in \mathbb{R}^{3 \times H_v \times W_v}$  and  $a \in \mathbb{R}^{H_a \times W_a}$  denote a video frame and a corresponding audio signal from the same

video clip, respectively. Here the raw 1D audio waveform has been converted into the 2D spectrogram by Short-Time Fourier Transform (STFT), and therefore we use 2D CNNs to extract deep semantic features of audio modality like vision. In practice, we employ the off-the-shelf VGG16 [47] for frame processing ( $\text{Enc}_v$ ) and the VGGish network [26] for spectrogram analysis ( $\text{Enc}_a$ ), similar to [29]. The output feature map of the final convolution layer of VGG16 is treated as the original visual feature  $f_v$ . We use layers before the final post-processing stage of VGGish to produce a high-level embedding as the original audio feature  $f_a$ . These feature extraction processes are formulated as:

$$f_v = \text{Enc}_v(I), \quad f_v \in \mathbb{R}^{c_v \times h \times w}, \quad (1)$$

$$f_a = \text{Enc}_a(a), \quad f_a \in \mathbb{R}^{c_a}. \quad (2)$$

### 3.2. Predictive Learning across Visual Views

This section details the vanilla SSPL, which contains the attention module to compute audio-visual representation from previously extracted features, and the self-supervised learning to guide model training via cross-view representation prediction.

Let  $I^1$  and  $I^2$  denote two randomly augmented views of the given image  $I$ . The two views are respectively fed into the visual CNN,  $\text{Enc}_v$ , to obtain spatial feature maps  $f_v^1$  and  $f_v^2$  as in Eq. (1). Besides, an audio feature vector  $f_a$  is derived from the audio signal  $a$  using Eq. (2). For simplicity, we use  $n \in \{1, 2\}$  to index different visual views.

**Attention Module (AM).** We adopt the normalized inner product (or cosine similarity) to measure the similarity between audio and visual features, as suggested by [5, 46]. Considering that  $f_a$  and  $f_v^n$  are from two heterogeneous modalities, we first transform  $f_a$  to be comparable with the visual feature via a nonlinear transformation,  $\tilde{f}_a = g(f_a) \in \mathbb{R}^{c_v}$ , and then perform the similarity measurement. Formally, for the spatial location  $(i, j)$  in visual feature map  $f_v^n$ , a similarity value is computed as follows:

$$S^n(i, j) = \frac{\langle \tilde{f}_a, f_v^n(\cdot, i, j) \rangle}{\|\tilde{f}_a\|_2 \|f_v^n(\cdot, i, j)\|_2}, \quad (i, j) \in [h] \times [w], \quad (3)$$

where  $f_v^n(\cdot, i, j) \in \mathbb{R}^{c_v}$ .

The similarity map  $S \in \mathbb{R}^{h \times w}$  plays two important roles in our method. On the one hand, it indicates the degree of correlation between each image location (after resized to image scale) and the given audio cues, and thus can serve as the sound localization map. On the other hand, it acts as an attention mechanism to weigh the original visual feature, resulting in the following audio-visual representation:

$$f_{av}^n(k) = \sum_{i,j} \tilde{S}^n(i, j) f_v^n(k, i, j), \quad k \in \{1, \dots, c_v\}, \quad (4)$$

$$\tilde{S}^n = \frac{S^n - \min(S^n)}{\max(S^n) - \min(S^n)}. \quad (5)$$

Here we scale the similarity values to  $[0, 1]$  by min-max normalization [35]. This operation makes different feature elements more distinguishable, and performs better compared with the sigmoid and softmax scaling functions [43, 46] (see Table 5 for empirical comparisons). Since the  $f_{av}^n \in \mathbb{R}^{c_v}$  selects and integrates visual features that are more related to audio ones, we treat it as a multi-modal representation to advance subsequent learning.

**Self-Supervised Learning.** The learning procedure aims to make the two audio-visual representations similar. Our hypothesis is that two visual scenes containing the same sounding objects should consistently correspond to the same audio cues in semantic level. We follow SimSiam [12] to achieve this goal in the audio-visual setting.

Formally, we feed  $f_{av}^n$  into a MLP head to obtain the projection of corresponding view,  $z^n = \text{MLP}(f_{av}^n)$ . Then a predictor head, denoted as  $\text{Pred}$ , takes as input  $z^1$  to predict  $z^2$  by minimizing the negative cosine similarity (NCS):

$$\mathcal{L}_{NCS}(z^1, z^2) = -\frac{\langle \text{Pred}(z^1), z^2 \rangle}{\|\text{Pred}(z^1)\|_2 \|z^2\|_2}. \quad (6)$$

To symmetrize the above loss, we also feed  $z^2$  into  $\text{Pred}$  to estimate  $z^1$ , leading to another loss term  $\mathcal{L}_{NCS}(z^2, z^1)$ . The total loss is therefore defined as:

$$\mathcal{L}_{SSPL} = \frac{1}{2} \mathcal{L}_{NCS}(z^1, z^2) + \frac{1}{2} \mathcal{L}_{NCS}(z^2, z^1). \quad (7)$$

However, as discussed in [12], directly minimizing the loss in Eq. (7) could easily induce representation collapse. To overcome this problem the stop-gradient (SG) operation is employed. That is, Eq. (6) is modified as  $\mathcal{L}_{NCS}(z^1, \text{SG}(z^2))$ , where  $z^2$  is viewed as a constant such that branch on  $I^2$  receives no gradient from  $z^2$  through this loss term. Similarly we have  $\mathcal{L}_{NCS}(z^2, \text{SG}(z^1))$ , and the form in Eq. (7) is implemented as:

$$\mathcal{L}_{SSPL} = \frac{1}{2} \mathcal{L}_{NCS}(z^1, \text{SG}(z^2)) + \frac{1}{2} \mathcal{L}_{NCS}(z^2, \text{SG}(z^1)). \quad (8)$$

Note that we follow the SimSiam framework for its simplicity in the use of only positive pairs without representation collapse. However, our predictive learning strategy can be combined with other self-supervised learning methods, such as BYOL [21], W-MSE [15], and BarlowTwins [55]. We leave the potential extension for future works.

### 3.3. Predictive Learning across Modalities

In this section, we propose the PCM for audio and visual feature alignment, and continuously improving the localization performance of SSPL. The key idea inherits the spirit of predictive coding (PC) in neuroscience [45, 48, 54], which simulates the mechanism of information processing in visual cortex. In brief, PC uses feedback connections from a



higher-level area to a lower-level one to convey predictions of lower-level neural activities; it employs feedforward connections to carry the errors between the actual activities and the predictions; and the brain dynamically updates representations so as to progressively reduce the prediction errors. In our PCM (Figure 3), we treat the visual feature as a type of prior knowledge to predict the audio feature in an *iterative* manner. In the following, at the heart of PCM, we give the representation update rules of feedback and forward processes, respectively. The detailed derivations can be found in supplement.

Denote by  $r_l(t)$ ,  $l \in \{1, \dots, L\}$ ,  $t \in \{0, \dots, T\}$  the representation of the  $l$ -th layer of PCM network at time step  $t$ , and by  $W_{l,l-1}$  the feedback connection weights from layer  $l$  to layer  $l-1$  (and vice versa for  $W_{l-1,l}$ ).

The **feedback process** updates representations through a mechanism of layer-wise prediction generation. Concretely, at  $l$ -th layer the prediction,  $p_l$ , of representation,  $r_l$ , is first derived using the above layer's representation,  $r_{l+1}$ . Then  $r_l$  is updated with its previous state and the prediction, *i.e.*, at time step  $t$  we have:

$$p_l(t) = (W_{l+1,l})^T r_{l+1}(t), \quad (9)$$

$$r_l(t) \leftarrow \phi((1 - b_l)r_l(t-1) + b_l p_l(t)), \quad (10)$$

where  $\phi$  is a nonlinear activation function and  $b_l$  serves as a positive scalar to balance two terms. The above update rules are executed from top layer  $L$  to bottom layer 1 in sequence, and by setting  $p_L(t) \equiv f_v^n$ , we in fact achieve further feature extraction from visual source.

In **feedforward process**, representations are again modulated based on prediction errors emerged at each layer. Specifically, the representation  $r_{l-1}$  and its prediction  $p_{l-1}$  are often unequal, resulting in a prediction error  $e_{l-1}$ . The error signal contains unpredictable components of  $r_{l-1}$ , and is forwarded to higher level to correct the representation  $r_l$ . This leads to complementary update rules:

$$e_{l-1}(t) = r_{l-1}(t) - p_{l-1}(t), \quad (11)$$

$$r_l(t) \leftarrow \phi(r_l(t) + a_l(W_{l-1,l})^T e_{l-1}(t)), \quad (12)$$

where  $r_0(t) \equiv f_a$  is the original audio feature,  $p_0(t) = \phi((W_{1,0})^T r_1(t))$  refers to the prediction of  $f_a$ , and  $a_l$  denotes a trade-off scalar like  $b_l$ .

PCM conducts the two distinct processes alternatively while all layers' representations are progressively refined so as to reduce the prediction error. Subsequently, we use a  $1 \times 1$  convolution to transform the top layer representation at last time step,  $r_L(T)$ , to a new visual feature,  $\tilde{f}_v^n$ , with the same dimension of  $f_v^n$ . Consequently the vanilla SSPL method in Sec. 3.2 can be enhanced by feeding  $\tilde{f}_v^n$ , instead of  $f_v^n$ , into the AM to compute audio-visual representation (*i.e.*, Eqs. (3) and (4)).

## 4. Experiments

### 4.1. Datasets and Evaluation Metrics

**SoundNet-Flickr** [4]. This dataset consists of more than 2 million videos from Flickr. We use a 3s audio clip around the middle frame of the whole audio, and the accompanied video frame to form an image-audio pair. Following [5, 46], we train models with two random subsets of 10k and 144k image-audio pairs, respectively, and perform evaluation on the 250 annotated pairs provided by [46]. Note that the location of the sound source in each test frame is given by 3 separate bounding boxes, each of which is obtained by a different annotator.

**VGG-Sound** [6] and **VGG-Sound Source** [5]. VGG-Sound dataset contains over 200k video clips that are divided into 300 sound categories. Similar to [5], we conduct training with 10k and 144k image-audio pairs randomly sampled from this dataset, respectively. For fair comparisons with recent works [1, 5, 46], we evaluate models on the VGG-Sound Source (VGG-SS) benchmark with 5k annotated image-audio pairs collected by [5]. Compared with SoundNet-Flickr benchmark that spans about 50 sounding object classes, VGG-SS has 220 classes and thus providing a more challenging scenario for sound localization task.

We focus on and reimplement two related methods, Attention [46] and HardWay [5] (SOTA on this task), which could be representatives of GLCL- and LLCL-based approaches, respectively. We denote our method without using PCM by SSPL (w/o PCM), and the version equipped with PCM by SSPL (w/ PCM). Additionally, we employ consensus Intersection over Union (cIoU) and Area Under Curve (AUC) as evaluation metrics, and report cIoU scores with threshold 0.5 in experiments, same as [5, 35, 43, 46].

### 4.2. Implementation Details

We use VGG16 [47] pretrained on ImageNet [14] and VGGish [26] pretrained on AudioSet [19] as visual and audio feature extractors, respectively. The visual input is an image of size  $256 \times 256 \times 3$ , on which we perform the data augmentation pipeline: random cropping with  $224 \times 224$  resizing and random horizontal flip. The raw 3s audio signal is re-sampled at 16kHz and further transformed into  $96 \times 64$  log-mel spectrograms as audio input, and the audio output feature  $f_a$  is a 128D vector. The non-linear audio feature transformation function  $g(\cdot)$  is instantiated with a simple two-layer network as in [46]: FC(512)-ReLU-FC(512). We closely follow SimSiam [12] to set the projection and prediction MLPs. For PCM, we mainly adopt Conv-MaxPool-GELU layers in the feedback pathway, and Upsample-DeConv-GELU layers in the feedforward counterpart. The weights of two feature extractors are kept frozen during training, and we optimize the rest of the model with AdamW [36]. We utilize the early stop-

Method	Training set	cIoU $\uparrow$	AUC $\uparrow$
Attention [46] <sup>†</sup> <sub>CVPR18</sub>	Flickr10k	0.442	0.461
DMC [27] <sub>CVPR19</sub>	Flickr10k	0.414	0.450
CAVL [29] <sub>arXiv20</sub>	Flickr10k	0.500	0.492
MSSL [43] <sub>ECCV20</sub>	Flickr10k	0.522	0.496
AVObject [1] <sub>ECCV20</sub>	Flickr10k	0.546	0.504
DSOL [28] <sub>NeurIPS20</sub>	Flickr10k	0.566	0.515
HardWay [5] <sup>†</sup> <sub>CVPR21</sub>	Flickr10k	0.615	0.535
ICL [35] <sub>CVPRW21</sub>	Flickr10k	0.710	0.580
SSPL (w/o PCM)	Flickr10k	0.671	0.556
SSPL (w/ PCM)	Flickr10k	<b>0.743</b>	<b>0.587</b>
Attention [46] <sub>CVPR18</sub>	Flickr144k	0.660	0.558
DMC [27] <sub>CVPR19</sub>	Flickr144k	0.671	0.568
HardWay [5] <sup>†</sup> <sub>CVPR21</sub>	Flickr144k	0.699	0.590
SSPL (w/o PCM)	Flickr144k	0.699	0.580
SSPL (w/ PCM)	Flickr144k	<b>0.759</b>	<b>0.610</b>
Attention [46] <sup>*</sup> <sub>CVPR18</sub>	VGG-Sound10k	0.522	0.502
HardWay [5] <sup>†</sup> <sub>CVPR21</sub>	VGG-Sound10k	0.647	0.560
SSPL (w/o PCM)	VGG-Sound10k	0.699	0.572
SSPL (w/ PCM)	VGG-Sound10k	<b>0.763</b>	<b>0.591</b>
HardWay [5] <sup>†</sup> <sub>CVPR21</sub>	VGG-Sound144k	0.723	<b>0.605</b>
HardWay [5] <sub>CVPR21</sub>	VGG-Sound Full	0.735	0.590
SSPL (w/o PCM)	VGG-Sound144k	0.739	0.602
SSPL (w/ PCM)	VGG-Sound144k	<b>0.767</b>	<b>0.605</b>

Table 1. **Quantitative localization results on SoundNet-Flickr test set.** “\*” denotes our reproduction, and “†” indicates *improved* reproduction vs. original papers (see supplement).

ping strategy to avoid overfitting in all cases. More setting details (*e.g.*, learning rate and batch size) are in supplement.

### 4.3. Comparisons with State-of-the-art Methods

We first compare SSPL with recent methods on the SoundNet-Flickr test set in Table 1. We observe that when trained by 10k Flickr samples, the vanilla SSPL (w/o PCM) performs favorably against the two competing methods, HardWay [5] and ICL [35], while the enhanced SSPL (w/ PCM) outperforms the previous best [35] by a large margin (0.710 vs. 0.743, around 5% improvement). In the Flickr144k training case, SSPL (w/ PCM) increases performance by 8.6% cIoU and 3.4% AUC compared to HardWay, establishing a new state-of-the-art on this benchmark. These results demonstrate that SSPL without relying on negatives is feasible and effective for sound localization. Following [5], we also train on VGG-Sound using respective 10k and 144k data pairs, which enables SSPLs to achieve the top two localization performance in both settings. As discussed in [5], the sounding objects are often visible in video clips from VGG-Sound, revealing that our method can benefit from the improved data quality. What’s more, the performance of SSPL is significantly boosted by PCM, especially in the 10k’s setting (11% improvement for Flickr10k and 9% for VGG-Sound10k). This illustrates the advantage of PCM for facilitating sound localization.

We further evaluate SSPL on the newly released VGG-

Method	Training set	cIoU $\uparrow$	AUC $\uparrow$
Attention [46] <sup>*</sup> <sub>CVPR18</sub>	VGG-Sound10k	0.160	0.283
HardWay [5] <sup>*</sup> <sub>CVPR21</sub>	VGG-Sound10k	0.277	0.349
SSPL (w/o PCM)	VGG-Sound10k	0.253	0.335
SSPL (w/ PCM)	VGG-Sound10k	<b>0.314</b>	<b>0.369</b>
Attention [46] <sup>*</sup> <sub>CVPR18</sub>	VGG-Sound144k	0.171	0.287
AVObject [1] <sub>ECCV20</sub>	VGG-Sound144k	0.297	0.357
HardWay [5] <sup>*</sup> <sub>CVPR21</sub>	VGG-Sound144k	0.319	0.370
SSPL (w/o PCM)	VGG-Sound144k	0.270	0.348
SSPL (w/ PCM)	VGG-Sound144k	<b>0.339</b>	<b>0.380</b>

Table 2. **Quantitative localization results on VGG-SS test set.**

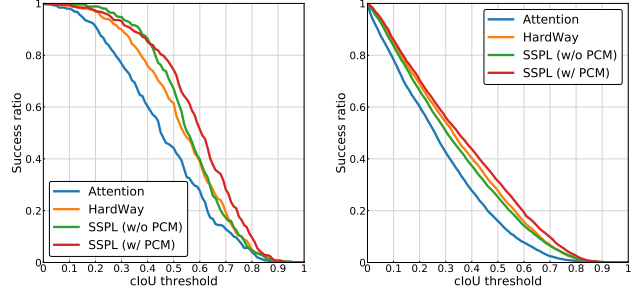


Figure 4. **Success ratio with varying cIoU thresholds.** Left: SoundNet-Flickr test set. Right: VGG-SS test set. Best viewed in color and by zooming in.

SS benchmark and report results in Table 2. Because in this challenging benchmark the sounding object categories are more diverse and the number of test samples is greater than those of SoundNet-Flickr [5], the performance of all methods drops severely compared with the results in Table 1. While SSPL (w/o PCM) still outperforms Attention by a large margin, it does not overtake HardWay. We attribute this to the limitation of vanilla SSPL on dealing with background noise (see Sec. 4.4 for an empirical comparison). However, by combining with feature alignment module, SSPL (w/ PCM) yields performance better than the state-of-the-art HardWay, especially by a substantial gap in the 10k’s scenario (0.277 vs. 0.314, over 13% gain). This verifies the superiority of the enhanced SSPL.

To address diverse demands for sound localization fineness, we compute cIoU scores with various thresholds as shown in Figure 4. The proposed method, SSPL (w/ PCM), again consistently surpasses the state-of-the-art (HardWay) under all thresholds.

### 4.4. Qualitative Analysis

We provide visualized localization results in Figure 5. We observe that Attention [46] is prone to overlook target objects (*e.g.*, the first and second rows in Figure 5a) and cover unrelated background details (*e.g.*, ground and sky). Since localization map also visualizes similarities between audio and visual features, the inaccurate localization indicates that Attention (random sampling of negatives) has the potential to misalign features. Although HardWay [5]

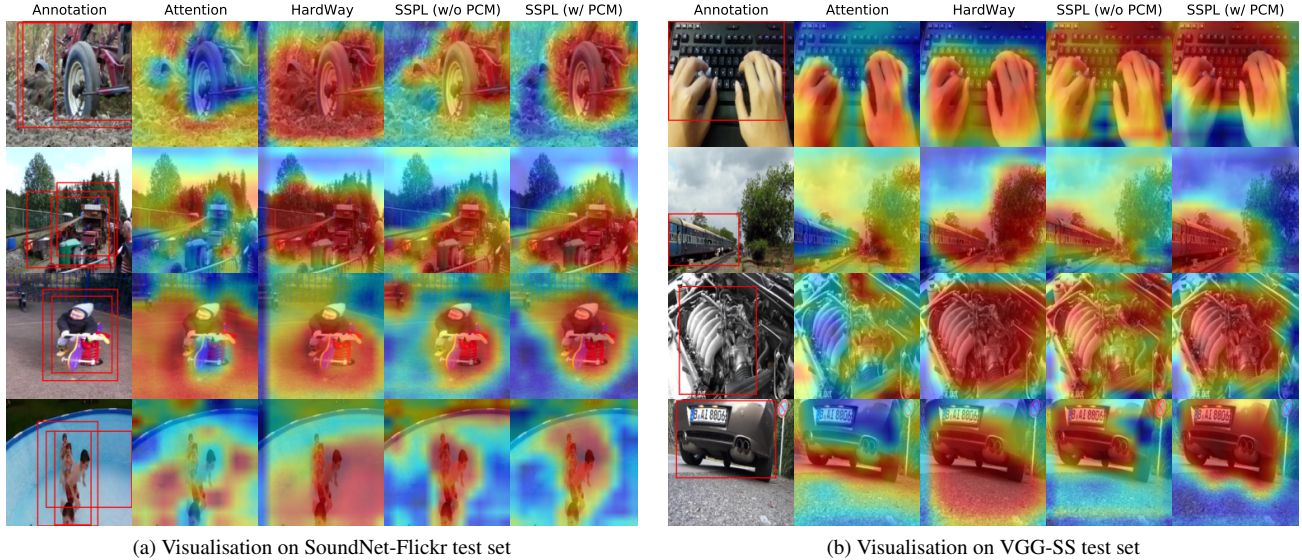


Figure 5. **Qualitative comparisons.** In each panel, the first column shows images accompanied with annotations, and remaining columns represent the predicted localization of sounding objects. Here the attention map or similarity map produced by different methods is visualized as the localization map. Note that for SoundNet-Flickr the bounding boxes are derived from multiple annotators.

	Pre-train	Stop-grad	$T$	cIoU $\uparrow$	AUC $\uparrow$
(a)				0.141	0.147
(b)	✓			0.382	0.432
(c)		✓		0.570	0.511
(d)	✓	✓		<b>0.671</b>	<b>0.556</b>
(e)	✓	✓	1	0.655	0.562
(f)	✓	✓	3	0.719	0.584
(g)	✓	✓	5	<b>0.743</b>	<b>0.587</b>

Table 3. **Ablation on training strategies.** “Pre-train” represents whether using the ImageNet-pretrained backbone to extract visual features, and  $T$  denotes the recursive cycles for iterative computing in PCM during training.

presents more centralized attention via hard negative mining, it easily underestimates (e.g., the first row in Figure 5b) or overestimates (e.g., the second row in Figure 5b) extents of sounding objects. This is probably because positive and negative regions in different images cannot be simply distinguished by the same thresholding parameters [5]. By contrast, our SSPL can cover the main region of interest, and the use of PCM further helps reduce the influence of background noise, leading to more accurate localization.

#### 4.5. Ablation Study

In this section, we delve deeper into SSPL by conducting extensive ablation studies. Unless otherwise specified, all experiments are performed on SoundNet-Flickr dataset.

**Training Strategy.** As discussed in prior art [12], a simple Siamese network without using negative samples can easily suffer from the problem of representation collapse. In this regard we evaluate key factors of SSPL that facilitate audio-

visual learning. In Table 3a we train the model from scratch while removing the stop-gradient operation, which indeed causes collapse in our practice. The variant with only pre-training strategy (Table 3b) improves performance because of the better parameter initialization, but it does not avoid collapsed solution yet. Adding stop-gradient alone during training (Table 3c) can obtain obvious gains, and the combination with pre-training (Table 3d) further boosts cIoU to 0.671, which is the default configuration of vanilla SSPL.

Based on above configuration, we perform additional ablation on the recursive cycles for representation updates in PCM. The performance slightly drops by 2% as conducting feedback and feedforward representation updates (Eqs. (9) to (12)) only once (Table 3e). This is because one computing step is not enough for PCM to reduce prediction errors between audio and visual features, and such non-negligible errors could degrade the subsequent learning. However, by increasing recursive cycles, SSPL can harvest significant performance improvements (nearly 10% in Table 3f and over 13% in Table 3g, respectively).

In summary, the results demonstrate that stop-gradient also works in our audio-visual setting to prevent collapse; and that both pre-training and PCM induce the model to learn effectively so as to promote localization accuracy.

**Augmentation.** We investigate the influence of various image augmentations on localization. As shown in Table 4, with the random crop baseline, our method can already achieve reasonable performance, indicating that object scales really matter in SSPL. However, except for horizontal flip (over 30% and 8% improvements on two datasets, respectively), randomly combining other augmen-



Augmentation	SoundNet-Flickr		VGG-SS	
	cIoU $\uparrow$	AUC $\uparrow$	cIoU $\uparrow$	AUC $\uparrow$
Crop (baseline)	0.514	0.499	<u>0.233</u>	0.324
+ Horizontal flip	<u>0.671</u>	<u>0.556</u>	<b>0.253</b>	<b>0.335</b>
+ Vertical flip	0.667	0.551	0.213	0.317
+ Translation	0.643	0.541	0.216	0.313
+ Rotation	0.639	0.543	0.227	<u>0.331</u>
+ Grayscale	0.610	0.535	0.226	0.318
+ Color jittering	<b>0.679</b>	<b>0.560</b>	0.232	0.328
+ Gaussian blur	0.619	0.533	0.204	0.299

Table 4. **Ablation on image augmentations.** Bold indicates the best and Underline the runner-up. Parameters used to generate different augmentations are provided in supplement.

Scaling method	cIoU $\uparrow$	AUC $\uparrow$
ReLU [43]	0.353	0.424
Sigmoid	0.647	0.547
Softmax [46]	0.667	0.554
ReLU + Softmax [46]	0.574	0.531
Min-Max Norm.	<b>0.671</b>	<b>0.556</b>

Table 5. **Ablation on scaling methods.**

tations with crop cannot obtain consistent gains. This is because compared with other combinations, the spatial augmentations (random crop + horizontal flip) are more suitable for the pretrained and frozen VGG [47] to extract semantic visual features. Since our work is inspired by SimSiam [12], we also adopt its data augmentation strategies in SSPL, but find no benefits in this setting. Therefore, in all experiments we take the spatial augmentations by default.

**Scaling Method in AM.** The similarity map takes values in  $[-1, 1]$  and is adapted to weigh visual features in AM. Here we study different methods that can scale similarity range into  $[0, 1]$ . ReLU is used in [43] to compact the similarities less than 0, but in our model enforcing those negative values to be equal produces worst results, as shown in Table 5. While sigmoid and softmax [46] boost performance by taking all different similarities into account, they shrink values into a proper subset of  $[0, 1]$ . The min-max normalization (Eq. (5)), by contrast, takes a step forward and separates minima and maxima to the largest extent, yielding best results among others. This reveals that the relative importance between spatial-wise visual features is more crucial than the feature value per se for sound localization task.

**Further Analysis of PCM.** We empirically clarify the remarkable ability of PCM to boost sound localization. Since PCM features an iterative computing procedure, we inspect the performance of SSPL (w/ PCM) with different iterations in Figure 6. We observe that the localization accuracy tends to increase given more iterative computations, especially at the initial three time steps. To understand why this is the case, we look into attention maps from some test samples, as shown in Figure 7. PCM infers different visual representations with varying time steps (1 through 5), which are

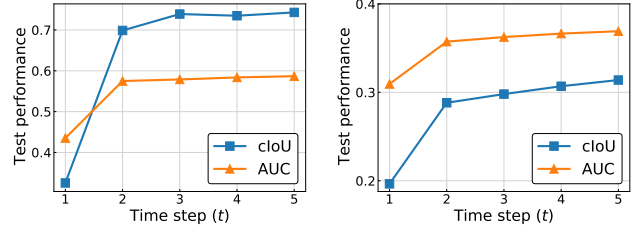


Figure 6. **Performance with PCM's iterations during testing.** Left: SoundNet-Flickr test set. Right: VGG-SS test set.

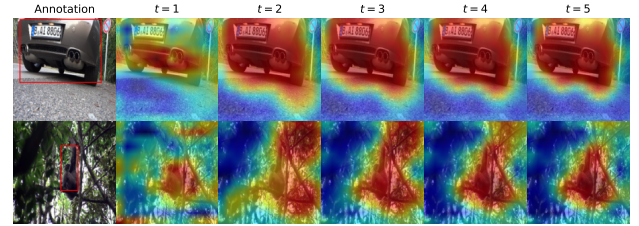


Figure 7. **Attention map with PCM's iterations during testing.** Illustrations are from VGG-SS test set.

further used by AM to yield different attention maps. Attention is less definitive (light red on sounding objects) and/or inaccurate (crimson on backgrounds) at early time steps. At later time steps, however, the model corrects itself to pay more definitive and accurate attention to the objects of interest. Adjusting attention in such a coarse-to-fine manner is particularly helpful to address ambiguous cases, where the object's appearance may be similar to backgrounds (e.g., the second row in Figure 7).

## 5. Conclusion and Future Works

In this work, we have developed a self-supervised audio-visual learning method, SSPL, that improves visual sound localization performance by explicit positive mining. A three-stream network, as well as its training strategy, was designed to explore correspondence between sound and video frame from the same video clip. We further proposed PCM to align audio and visual features via cross-modal feature prediction, which boosts localization accuracy significantly. Our approach shows promising performance on sound localization task, especially achieving the new state-of-the-art on SoundNet-Flickr benchmark.

While SSPL excels at single sound source localization, it is not applicable to localize multiple sound sources in unconstrained videos [43], which is still a challenge for the community. A potential solution is to develop weakly- or semi-supervised methods. We leave it for future works.

**Acknowledgements.** This work was supported in part by the Major Project for New Generation of AI (No. 2018AAA0100400), in part by the National Natural Science Foundation of China (Nos. 61836014, U21B2042, 62072457, 62006231, and 61976174), and in part by the Project funded by China Postdoctoral Science Foundation (No. 2021M703489).



## References

- [1] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In *ECCV*, pages 208–224, 2020. 1, 2, 3, 5, 6
- [2] Humam Alwassel, Dhruv Mahajan, Bruno Korbar, Lorenzo Torresani, Bernard Ghanem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. In *NeurIPS*, pages 9758–9770, 2020. 3
- [3] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *ICCV*, pages 609–617, 2017. 3
- [4] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. SoundNet: Learning sound representations from unlabeled video. In *NeurIPS*, pages 892–900, 2016. 3, 5
- [5] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *CVPR*, pages 16867–16876, 2021. 1, 2, 3, 4, 5, 6, 7, 11, 13, 14
- [6] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. VGGSound: A large-scale audio-visual dataset. In *ICASSP*, pages 721–725, 2020. 5
- [7] Peihao Chen, Deng Huang, Dongliang He, Xiang Long, Runhao Zeng, Shilei Wen, Mingkui Tan, and Chuang Gan. RSPNet: Relative speed perception for unsupervised video representation learning. In *AAAI*, pages 1045–1053, 2021. 2
- [8] Peihao Chen, Yang Zhang, Mingkui Tan, Hongdong Xiao, Deng Huang, and Chuang Gan. Generating visually aligned sound from videos. *IEEE Transactions on Image Processing*, 29:8292–8302, 2020. 3
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607, 2020. 2
- [10] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. In *NeurIPS*, pages 22243–22255, 2020. 2
- [11] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 2
- [12] Xinlei Chen and Kaiming He. Exploring simple Siamese representation learning. In *CVPR*, pages 15750–15758, 2021. 2, 4, 5, 7, 8, 13
- [13] Yanbei Chen, Yongqin Xian, A Sophia Koepke, Ying Shan, and Zeynep Akata. Distilling audio-visual knowledge by compositional contrastive learning. In *CVPR*, pages 7016–7025, 2021. 3
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. 5
- [15] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. In *ICML*, pages 3015–3024, 2021. 2, 4
- [16] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *CVPR*, pages 3299–3309, 2021. 2
- [17] John W Fisher III, Trevor Darrell, William T Freeman, and Paul A Viola. Learning joint statistical models for audio-visual fusion and segregation. In *NeurIPS*, pages 772–778, 2000. 1, 3
- [18] Chuang Gan, Hang Zhao, Peihao Chen, David Cox, and Antonio Torralba. Self-supervised moving vehicle tracking with stereo sound. In *ICCV*, pages 7053–7062, 2019. 3
- [19] Jort F Gemmeke, Daniel P W Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio Set: An ontology and human-labeled dataset for audio events. In *ICASSP*, pages 776–780, 2017. 5
- [20] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Un-supervised representation learning by predicting image rotations. In *ICLR*, 2018. 13
- [21] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, pages 21271–21284, 2020. 2, 4
- [22] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 2
- [23] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S M Ali Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. In *ICML*, pages 4182–4192, 2020. 2
- [24] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2016. 12
- [25] John Hershey and Javier Movellan. Audio-Vision: Using audio-visual synchrony to locate sounds. In *NeurIPS*, pages 813–819, 1999. 1, 3
- [26] Shawn Hershey, Sourish Chaudhuri, Daniel P W Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, Malcolm Slaney, Ron J Weiss, and Kevin Wilson. CNN architectures for large-scale audio classification. In *ICASSP*, pages 131–135, 2017. 4, 5
- [27] Di Hu, Feiping Nie, and Xuelong Li. Deep multimodal clustering for unsupervised audiovisual learning. In *CVPR*, pages 9248–9257, 2019. 1, 3, 6
- [28] Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou. Discriminative sounding objects localization via self-supervised audiovisual matching. In *NeurIPS*, pages 10077–10087, 2020. 1, 6
- [29] Di Hu, Zheng Wang, Haoyi Xiong, Dong Wang, Feiping Nie, and Dejing Dou. Curriculum audiovisual learning. *arXiv preprint arXiv:2001.09414*, 2020. 3, 4, 6
- [30] Hamid Izadinia, Imran Saleemi, and Mubarak Shah. Multimodal analysis for identification and segmentation of moving-sounding objects. *IEEE Transactions on Multimedia*, 15(2):378–390, 2012. 3
- [31] Minsong Ki, Youngjung Uh, Junsuk Choe, and Hyeran Byun. Contrastive attention maps for self-supervised co-localization. In *ICCV*, pages 2803–2812, 2021. 13

- [32] Einat Kidron, Yoav Y Schechner, and Michael Elad. Pixels that sound. In *CVPR*, pages 88–95, 2005. 3
- [33] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, pages 7774–7785, 2018. 3
- [34] César Laurent, Gabriel Pereyra, Philémon Brakel, Ying Zhang, and Yoshua Bengio. Batch normalized recurrent neural networks. In *ICASSP*, pages 2657–2661, 2016. 12
- [35] Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, and Ming-Hsuan Yang. Unsupervised sound localization via iterative contrastive learning. In *CVPR Workshop*, 2021. 2, 3, 4, 5, 6
- [36] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 5, 12
- [37] Pedro Morgado, Yi Li, and Nuno Vasconcelos. Learning representations from audio-visual spatial alignment. In *NeurIPS*, pages 4733–4744, 2020. 2, 3
- [38] Pedro Morgado, Nuno Vasconcelos, and Ishan Misra. Audio-visual instance discrimination with cross-modal agreement. In *CVPR*, pages 12475–12486, 2021. 3
- [39] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010. 12
- [40] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2
- [41] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, pages 631–648, 2018. 1, 3
- [42] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *ECCV*, pages 801–816, 2016. 3
- [43] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, pages 292–308, 2020. 1, 2, 3, 4, 5, 6, 8
- [44] Rui Qian, Tianjian Meng, Boqing Gong, Ming-Hsuan Yang, Huisheng Wang, Serge Belongie, and Yin Cui. Spatiotemporal contrastive video representation learning. In *CVPR*, pages 6964–6974, 2021. 2
- [45] Rajesh P N Rao and Dana H Ballard. Predictive coding in the visual cortex: A functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87, 1999. 4
- [46] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In *CVPR*, pages 4358–4366, 2018. 1, 2, 3, 4, 5, 6, 8, 11, 13, 14
- [47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 4, 5, 8
- [48] Zengjie Song, Jianshe Zhang, Guang Shi, and Junmin Liu. Fast inference predictive coding: A novel model for constructing deep neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 30(4):1150–1165, 2018. 4
- [49] Michael W Spratling. A review of predictive coding algorithms. *Brain and Cognition*, 112:92–97, 2017. 11
- [50] Kun Su, Xiulong Liu, and Eli Shlizerman. How does it sound? Generation of rhythmic soundtracks for human movement videos. In *NeurIPS*, 2021. 3
- [51] Yapeng Tian, Di Hu, and Chenliang Xu. Cyclic co-learning of sounding object visual grounding and sound separation. In *CVPR*, pages 2745–2754, 2021. 1, 2, 3
- [52] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, pages 776–794, 2020. 2
- [53] Francisco Rivera Valverde, Juana Valeria Hurtado, and Abhinav Valada. There is more than meets the eye: Self-supervised multi-object detection and tracking with sound by distilling multimodal knowledge. In *CVPR*, pages 11612–11621, 2021. 3
- [54] Haiguang Wen, Kuan Han, Junxing Shi, Yizhen Zhang, Eugenio Culurciello, and Zhongming Liu. Deep predictive coding network for object recognition. In *ICML*, pages 5266–5275, 2018. 4, 11, 12
- [55] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow Twins: Self-supervised learning via redundancy reduction. In *ICML*, pages 12310–12320, 2021. 2, 4
- [56] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, pages 570–586, 2018. 2, 15

## Supplementary Material

This supplementary material contains four parts:

- Section A presents full derivations to formulate the representation update rules of predictive coding module (PCM).
- Section B provides more details on our implementation.
- Section C compares the localization performance of Attention [46] and HardWay [5] in original papers with our reproductions.
- Section D gives additional ablation and visualisation results.

### A. Full Formulation of PCM

The PCM, proposed for audio and visual feature alignment, plays an important role in improving sound localization performance of SSPL. As shown in Figure S1, the key idea underlying PCM consists of three parts: (1) a feedback process (solid line) updates representations with the top-down predictions that originate from the visual feature; (2) a feedforward process (dashed line) also updates representations but with the bottom-up prediction errors that evolve from the audio feature; (3) a recursive modulation mechanism works to conduct the two processes alternatively. In the following, we first formulate the optimization objective of PCM, and then derive the representation update rules of the two processes, respectively, which are followed by a brief summary and a formal algorithm. *Note that for applications of PCM, we only need to explicitly update representations according to the rules given in Eqs. (S10) to (S13), without performing derivations again.*

Denote by  $f_a$  the audio feature, by  $f_v$  the visual feature, by  $r_l(t), l \in \{1, \dots, L\}, t \in \{0, \dots, T\}$  the representation of the  $l$ -th layer of PCM network at time step  $t$ , and by  $W_{l,l-1}$  the feedback connection weights from layer  $l$  to layer  $l-1$  (and vice versa for  $W_{l-1,l}$ ).

**Optimization Objective.** At layer  $l$ , PCM minimizes the following compound loss:

$$\mathcal{L}_{PCM}^l = \frac{\alpha_l}{2} \underbrace{\|r_{l-1} - \mathcal{G}((W_{l,l-1})^T r_l)\|_2^2}_{\mathcal{L}_1^l} + \frac{\beta_l}{2} \underbrace{\|r_l - p_l\|_2^2}_{\mathcal{L}_2^l}, \quad (\text{S1})$$

where the function  $\mathcal{G}$  corresponds to a generative process,  $\alpha_l$  and  $\beta_l$  are scalars that control the weights of the two loss terms  $\mathcal{L}_1^l$  and  $\mathcal{L}_2^l$ , and  $p_l = \mathcal{G}((W_{l+1,l})^T r_{l+1})$  is the prediction of  $r_l$ .

Given the lower-level representation  $r_{l-1}$  and the top-down prediction  $p_l$ , our goal is to estimate  $r_l$  so as to decrease the loss in Eq. (S1). Minimizing  $\mathcal{L}_1^l$  w.r.t.  $r_l$  leads to

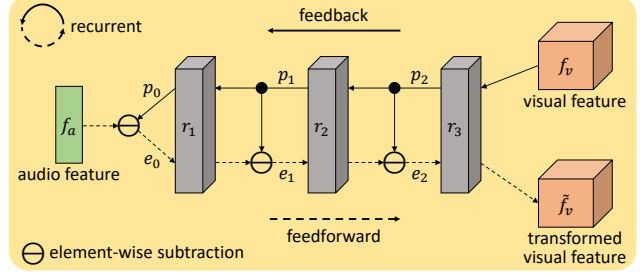


Figure S1. **Overview of predictive coding module (PCM).** For simplicity we only show a 3-layer version.

the representation that can be used to predict the *lower level* of representation  $r_{l-1}$ , while minimizing  $\mathcal{L}_2^l$  w.r.t.  $r_l$  yields the representation that approximates the prediction signal  $p_l$  coming from a *higher level*. Therefore, the representation  $r_l$  associates lower- and higher-level information by reducing two prediction errors in  $\mathcal{L}_1^l$  and  $\mathcal{L}_2^l$ . Minimizing losses at all layers can implicitly drive predictions at different levels to be mutually consistent [49].

**Feedback Process.** This process acts to update representations based on predictions from higher levels. Following [54], we set  $\mathcal{G}(x) = x$ , and then employ gradient descent to minimize  $\mathcal{L}_2^l$  w.r.t.  $r_l$ , resulting in update rules:

$$p_l(t) = (W_{l+1,l})^T r_{l+1}(t), \quad (\text{S2})$$

$$\frac{\partial \mathcal{L}_2^l}{\partial r_l(t)} = 2(r_l(t) - p_l(t)), \quad (\text{S3})$$

$$\begin{aligned} r_l(t+1) &= r_l(t) - \eta_l \frac{\beta_l}{2} \frac{\partial \mathcal{L}_2^l}{\partial r_l(t)} \\ &= (1 - \eta_l \beta_l) r_l(t) + \eta_l \beta_l p_l(t), \end{aligned} \quad (\text{S4})$$

where  $\eta_l$  is a non-negative scalar governing learning. For simplicity, let  $b_l = \eta_l \beta_l$ , and then Eq. (S4) is rewritten as follows:

$$r_l(t+1) = (1 - b_l) r_l(t) + b_l p_l(t). \quad (\text{S5})$$

PCM carries out the feedback updating from top layer  $L$  to bottom layer 1, where the prediction of  $r_L(t)$  at top layer is set as the visual feature, i.e.,  $p_L(t) \equiv f_v$ .

**Feedforward Process.** This process works to update representations by using prediction errors from lower levels. For layer  $l$ , the lower-level prediction error  $e_{l-1}$  is the difference between  $r_{l-1}$  and  $p_{l-1}$ . We use gradient decent to



minimize  $\mathcal{L}_1^l$  w.r.t.  $r_l$ , leading to the following update rules:

$$e_{l-1}(t) = r_{l-1}(t) - p_{l-1}(t), \quad (\text{S6})$$

$$\frac{\partial \mathcal{L}_1^l}{\partial r_l(t)} = -2W_{l,l-1}e_{l-1}(t), \quad (\text{S7})$$

$$\begin{aligned} r_l(t+1) &= r_l(t) - \kappa_l \frac{\alpha_l}{2} \frac{\partial \mathcal{L}_1^l}{\partial r_l(t)} \\ &= r_l(t) + \kappa_l \alpha_l W_{l,l-1} e_{l-1}(t), \end{aligned} \quad (\text{S8})$$

where  $\kappa_l$  is a non-negative scalar like  $\eta_l$ . We also set  $a_l = \kappa_l \alpha_l$  for simplicity. Similar to [54], we replace the feedback connection weights  $W_{l,l-1}$  in Eq. (S8) with the transposed feedforward connection weights  $(W_{l-1,l})^T$ , and thus can endow PCM with more degrees of freedom to learn. Consequently the update rule in Eq. (S8) can be rewritten as a feedforward operation:

$$r_l(t+1) = r_l(t) + a_l (W_{l-1,l})^T e_{l-1}(t). \quad (\text{S9})$$

In this process, PCM updates representations from bottom layer 1 to top layer  $L$ , where we let  $r_0(t) \equiv f_a$  and  $p_0(t) = (W_{1,0})^T r_1(t)$ .

**Summary and Algorithm.** So far we formulate PCM with the simple linear activation functions. To introduce non-linearity into PCM, a nonlinear activation function  $\phi$  (e.g., ReLU [39] used in [54] or GELU [24] used in this work) is applied to the above update Eqs. (S5) and (S9). By taking the recursive computing into account, we summarize the two processes as follows.

*Nonlinear feedback process* ( $l = L, L-1, \dots, 1$ ):

$$p_l(t) = (W_{l+1,l})^T r_{l+1}(t), \quad (\text{S10})$$

$$r_l(t) \leftarrow \phi((1 - b_l)r_l(t-1) + b_l p_l(t)). \quad (\text{S11})$$

*Nonlinear feedforward process* ( $l = 1, 2, \dots, L$ ):

$$e_{l-1}(t) = r_{l-1}(t) - p_{l-1}(t), \quad (\text{S12})$$

$$r_l(t) \leftarrow \phi(r_l(t) + a_l (W_{l-1,l})^T e_{l-1}(t)). \quad (\text{S13})$$

The two processes are conducted alternatively such that all representations in PCM are refined progressively. Finally, we transform the top layer representation at last time step,  $r_L(T)$ , to a new visual feature,  $\tilde{f}_v$ , with dimension the same as  $f_v$  by a  $1 \times 1$  convolution. The representation learning of SSPL can proceed based on this  $\tilde{f}_v$ , instead of  $f_v$  as used in the vanilla SSPL. We present main computing steps of PCM in Algorithm S1.

## B. Implementation Details

### B.1. Architecture of PCM

For the feedback process of PCM, we use convolution layers (kernel\_size = 3, stride = 1, padding =

---

### Algorithm S1 Update Representations in PCM

---

**Input:**  $f_v$  and  $f_a$

**Output:**  $\tilde{f}_v$

---

```

1: for  $t = 0$  to  $T$  do
2:   if  $t = 0$  then
3:     | initialize representations
4:   end if
5:   for  $l = L$  to  $1$  do ▷feedback process
6:     if  $l = L$  then
7:       |  $p_l(t) = f_v$ 
8:     else
9:       | compute prediction  $p_l(t)$ : Eq. (S10)
10:    end if
11:    update representation  $r_l(t)$ : Eq. (S11)
12:  end for
13:  for  $l = 1$  to  $L$  do ▷feedforward process
14:    if  $l = 1$  then
15:      |  $e_{l-1}(t) = f_a - \phi((W_{l,l-1})^T r_l(t))$ 
16:    else
17:      | obtain prediction error  $e_{l-1}(t)$ : Eq. (S12)
18:    end if
19:    update representation  $r_l(t)$ : Eq. (S13)
20:  end for
21: end for
22:  $\tilde{f}_v = \text{Conv}_{1 \times 1}(r_L(T))$  ▷transformed feature

```

---

1) followed by max pooling operation to reduce the spatial dimensionality of feature maps, while using  $1 \times 1$  convolutions to decrease the number of channels. As for the feed-forward process, the transposed convolutions (a.k.a. deconvolutions) are utilized and feature maps are upsampled by the “bilinear” upsampling algorithm, provided in PyTorch. Besides, the number of convolution layers is  $L = 3$ . From top layer  $L$  to bottom layer 1, the number of filters within each layer is 512, 512, and 128, respectively. The transposed convolution layers have the same setting. Moreover, we use GELU [24] as the nonlinear activation function for both processes. To stabilize and accelerate training, we adopt the batch normalization [34] before every non-linearity at each layer and at each time step, except the prediction of audio feature at bottom layer.

### B.2. Training Details for SSPL

The AdamW [36] optimizer is employed to train our model, where we set  $(\beta_1, \beta_2) = (0.9, 0.999)$  and set weight decay to  $10^{-4}$ . In practice, we find that better performance could be achieved if the learning rate for projection and prediction MLPs is greater than that for remaining model parts. We show detailed learning rate settings in Table S1. During training, there are 256 image-audio pairs in each minibatch, which are distributed in parallel on 2 or 4 NVIDIA GeForce GTX 1080 Ti GPUs.

Training set	SSPL (w/o PCM)		SSPL (w/ PCM)	
	$lr_1$	$lr_2$	$lr_1$	$lr_2$
SoundNet-Flickr	$2 \cdot 10^{-3}$	$5 \cdot 10^{-4}$	$5 \cdot 10^{-5}$	$2 \cdot 10^{-5}$
VGG-Sound	$1 \cdot 10^{-2}$	$5 \cdot 10^{-3}$	$5 \cdot 10^{-5}$	$2 \cdot 10^{-5}$

Table S1. **Learning rate settings.**  $lr_1$  denotes the learning rate for projection and prediction MLPs,  $lr_2$  for remaining model parts.

Augmentation	Parameter
Crop	$p = 1$
	output size of <code>Resize</code> = <code>int(224 × 1.1)</code>
	interpolation method of <code>Resize</code> = <code>BICUBIC</code>
	crop size = 224
Horizontal flip	$p = 0.5$
Vertical flip	$p = 0.5$
Translation	$p = 1.0$
	maximum absolute fraction = (0.2, 0.2)
Rotation	$p = 1.0$
	angle $\in \{0, 90, 180, 270\}$
Grayscale	$p = 0.2$
Color jittering	$p = 0.8$
	maximum brightness adjustment = 0.4
	maximum contrast adjustment = 0.4
	maximum saturation adjustment = 0.4
Gaussian blur	$p = 0.5$
	$\sigma \in [0.1, 2.0]$

Table S2. **Parameters used to generate image augmentations.**  $p$  denotes the probability that the corresponding operation will be performed.

### B.3. Image Augmentations for SSPL

As shown in Table S2, a total of 8 image augmentations are considered in our method. We follow HardWay [5] to select and set the first two augmentations: cropping with  $224 \times 224$  resizing and horizontal flip. Then, we verify the effectiveness of other three spatial augmentations that are widely used in self-supervised visual representation learning [20, 31], *i.e.*, vertical flip, translation, and rotation. Additionally, since our work draws inspiration from SimSiam [12], we also take into account its augmentation strategies: grayscale, color jittering, and Gaussian blur, while keeping their settings the same as SimSiam.

## C. Reproducing Related Methods

The quantitative comparisons in the main text are partially based on our reproductions of two related methods: Attention [46] and HardWay [5]. We reimplement them as faithfully as possible by following each corresponding paper. As show in Table S3, we are able to improve these two methods on SoundNet-Flickr by small and straightforward modifications. Specifically, we use Crop and Horizontal flip

Method	Source	Training set	cIoU $\uparrow$	AUC $\uparrow$
Attention [46]	O	Flickr10k	0.436	0.449
	R	Flickr10k	0.442	0.461
HardWay [5]	O	Flickr10k	0.582	0.525
	R	Flickr10k	0.615	0.535
HardWay [5]	O	Flickr144k	0.699	0.573
	R	Flickr144k	0.699	0.590
Attention [46]	O	VGG-Sound10k	-	-
	R	VGG-Sound10k	0.522	0.502
HardWay [5]	O	VGG-Sound10k	0.618	0.536
	R	VGG-Sound10k	0.647	0.560
HardWay [5]	O	VGG-Sound144k	0.719	0.582
	R	VGG-Sound144k	0.723	0.605

Table S3. **Our reproductions vs. original papers’ results on SoundNet-Flickr test set.** “O” denotes results from original papers and “R” our reproductions.

Method	Source	Training set	cIoU $\uparrow$	AUC $\uparrow$
Attention [46]	*	VGG-Sound144k	0.185	0.302
	R	VGG-Sound144k	0.171	0.287
HardWay [5]	O	VGG-Sound144k	0.344	0.382
	R	VGG-Sound144k	0.319	0.370

Table S4. **Our reproductions vs. original papers’ results on VGG-SS test set.** “\*” denotes results obtained from [5].

given in Table S2 to spatially augment images for Attention (*vs.* originally  $320 \times 320$  resizing), same as HardWay and our method. We also fine tune the learning rate and weight decay for these two competitors in order to achieve their best performance.

Table S4 compares our reproductions with original papers’ results on VGG-SS. In this case our reproductions are slightly lower than the original counterparts. Note that we had tried to adjust various hyper-parameters for HardWay training (*e.g.*, learning rate, weight decay, batch size, and number of training epochs) for multiple times, but better performance than reproductions shown in Table S4 were not achieved. We contribute the performance discrepancy to the updated data in this benchmark. On the one hand, Chen *et al.* [5] originally provides 5158 YouTube video IDs for testing, and users need to download, extract, and pre-process the designated audio and visual sources themselves. However, 466 videos (9%) are not available (removed or prohibited download) at the time of conducting our experiments, leading to 4692 image-audio pairs for final testing. On the other hand, as clarified by the authors on the official project page<sup>1</sup>, some bounding box annotations are updated recently and consequently it could cause a 2%-3% difference on performance. Based on these two aspects, we think that our reproductions are reasonable.

<sup>1</sup><https://github.com/hche11/Localizing-Visual-Sounds-the-Hard-Way>.

Fusion method	Cat	$\otimes$	$\oplus$	AM (ours)
cIoU $\uparrow$	0.285	0.538	0.647	<b>0.671</b>
AUC $\uparrow$	0.414	0.512	0.540	<b>0.556</b>

Table S5. **Ablation on feature fusion methods.** We use different fusion methods in SSPL (w/o PCM), and train models on SoundNet-Flickr10k while evaluating on the standard benchmark.

$T$	1	3	5	6	7	8
cIoU $\uparrow$	0.655	0.719	0.743	0.743	0.759	0.747
AUC $\uparrow$	0.562	0.584	0.587	0.595	0.595	0.590
GFLOPs $\downarrow$	38.3	43.0	47.6	49.9	52.2	54.5

Table S6. **Influence of recursive cycles  $T$  in PCM.** All models are trained on SoundNet-Flickr10k and evaluated on the standard benchmark.

## D. Additional Results

### D.1. Ablation on Feature Fusion Methods

In SSPL, visual and audio features are fused by the attention mechanism to compute audio-visual representation. Here we compare other three feature fusion methods, *i.e.*, concatenation (Cat), multiplication ( $\otimes$ ), and addition ( $\oplus$ ), with our attention module (AM). We can see from Table S5 that our AM outperforms others by a large margin. This verifies efficacy of the attention-based feature interaction.

### D.2. Balance between Performance and Complexity of PCM

In Table S6, we quantitatively compare performance and time complexity of SSPL with varying recursive cycles  $T$ . We find that more recursive cycles cannot always bring gains as performance tends to be saturated when  $T > 5$ . Additionally, compared with SSPL (w/o PCM) that occupies 35.9 GFLOPs, SSPL (w/ PCM) conducts more operations with increasing  $T$ . As shown in Table S6, PCM takes, on average, 2.3 GFLOPs to complete one iteration. To balance between performance and time complexity, we set  $T = 5$  during training.

### D.3. Effect of False Negatives on Localization

As discussed in the main text, learning with false negatives can induce ambiguity in localization results. In this section, we give more examples to empirically illustrate this effect. As shown in Figure S2, when the false negatives are allowed to take part in contrastive learning, sounding objects are easily ignored in final localization maps (method A). Although learning with true positive and negative samples harvests accurate localization, it requires class label to direct negative sampling (method B). By contrast, our method is able to obtain consistent localization among different image-audio pairs, without using negatives and labels at all (method C).

### D.4. Additional Qualitative Comparisons

In Figure S3, we illustrate more localization examples from Attention [46], HardWay [5], and our method SSPL on two standard benchmarks: SoundNet-Flickr and VGG-SS. Qualitative evaluation results show that our method can localize the full extent of sounding objects, especially for SSPL (w/ PCM) that yields more accurate localization by ignoring background noise.



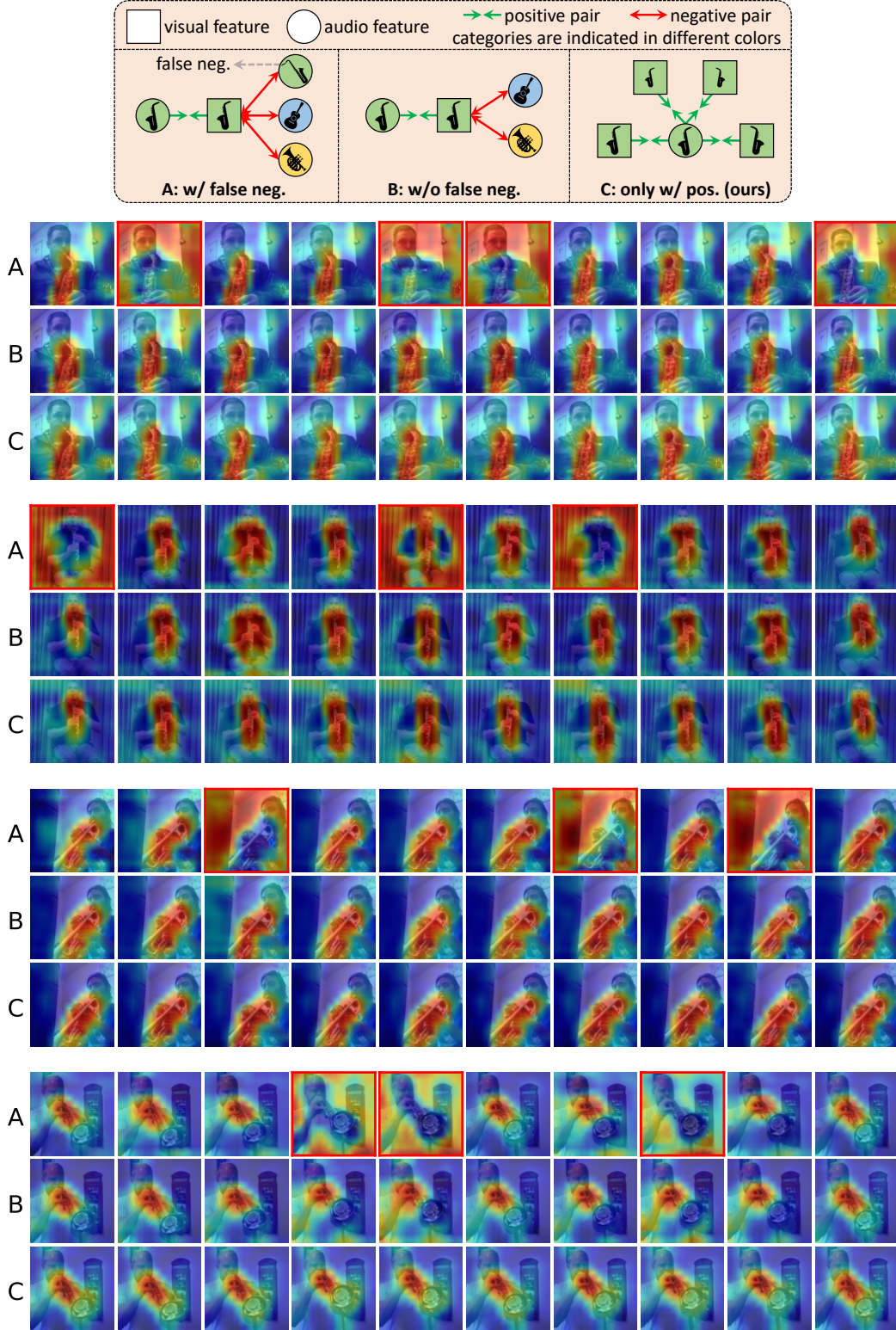
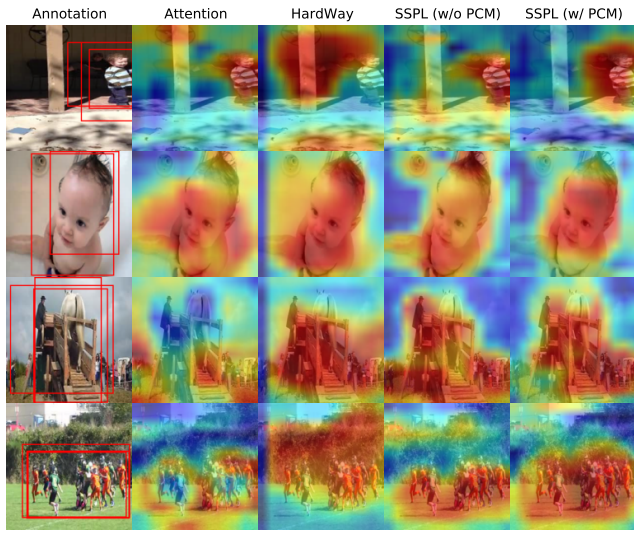
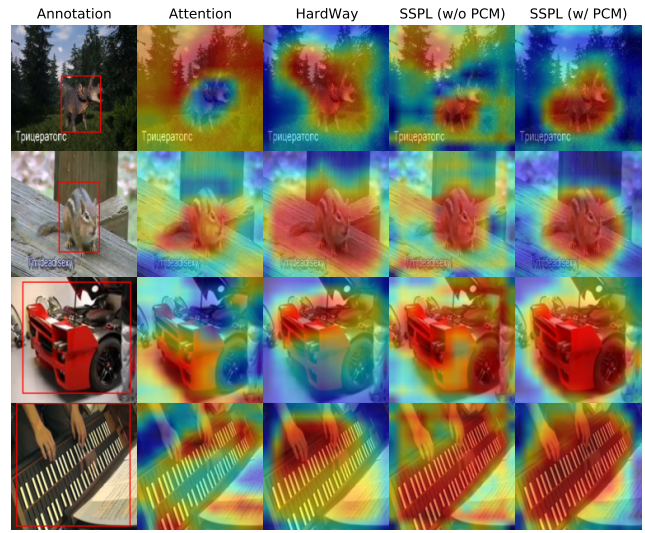


Figure S2. **Visualisation of the effect of false negatives on sound localization.** **A** denotes the training strategy that uses both true and false negatives to perform contrastive learning; **B** indicates the method where false negatives do not take part in contrastive learning, but *requiring class label to direct negative sampling*; **C** corresponds to our self-supervised method that only explores audio-visual positive pairs during learning. Here the false negatives are other videos' sounds that belong to the same category as the positive one. The images marked with red rectangle illustrate ambiguous localization results of method **A**. Models are trained on MUSIC [56].



(a) Visualisation on SoundNet-Flickr test set



(b) Visualisation on VGG-SS test set

Figure S3. **Qualitative comparisons.** In each panel, the first column shows images accompanied with annotations, and remaining columns represent the predicted localization of sounding objects. Here the attention map or similarity map produced by different methods is visualized as the localization map. Note that for SoundNet-Flickr the bounding boxes are derived from multiple annotators.