# Exploiting Transformation Invariance and Equivariance for Self-supervised Sound Localisation

Jinxiang Liu
Cooperative Medianet Innovation Center
Shanghai Jiao Tong University
China
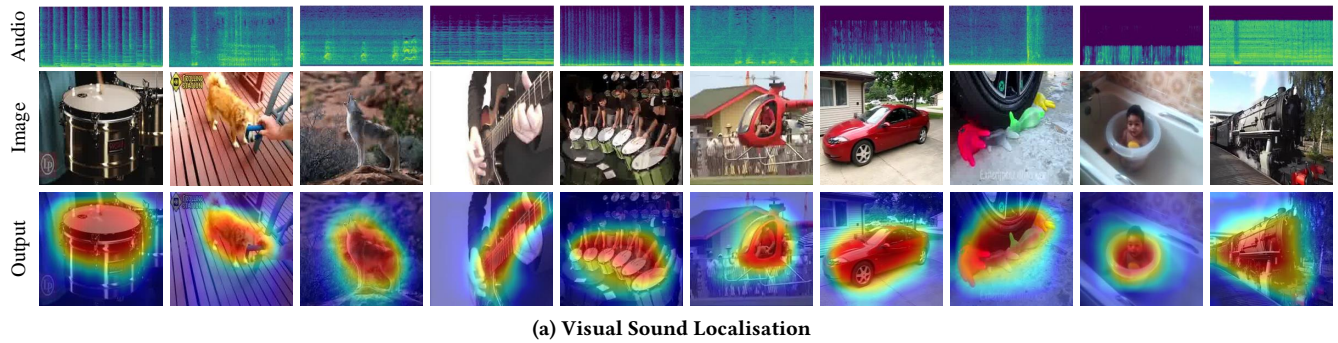jinxliu@sjtu.edu.cn

Chen Ju
Cooperative Medianet Innovation Center
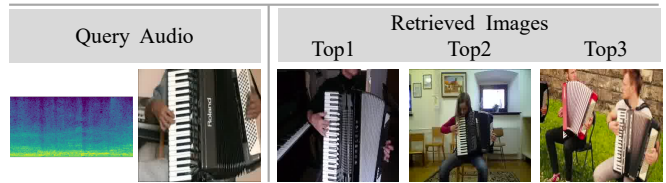Shanghai Jiao Tong University
China
ju_chen@sjtu.edu.cn

Weidi Xie*
[1]Cooperative Medianet Innovation Center
Shanghai Jiao Tong University
China
[2]Shanghai AI Laboratory
China
weidi@sjtu.edu.cn

Ya Zhang*
[1]Cooperative Medianet Innovation Center
Shanghai Jiao Tong University
China
[2]Shanghai AI Laboratory
China
ya_zhang@sjtu.edu.cn

**(a) Visual Sound Localisation**



**(b) Audio Retrieval**



**(c) Cross-modal Retrieval**

**Figure 1: Visualized results of our framework which learns powerful multi-modal representations for various applications. (a) visual sound localisation: highlights the salient object by its emitted sound. (b) audio retrieval: discovers semantic-identical audios to the query audio. (c) cross-modal retrieval: discovers semantic-identical images to the query audio.**

## ABSTRACT

We present a simple yet effective self-supervised framework for audio-visual representation learning, to localize the sound source in videos. To understand what enables to learn useful representations, we systematically investigate the effects of data augmentations, and reveal that (1) composition of data augmentations plays a critical role, *i.e.* explicitly encouraging the audio-visual representations to be invariant to various transformations (*transformation invariance*); (2) enforcing geometric consistency substantially improves the quality of learned representations, *i.e.* the detected sound source should follow the same transformation applied on input video frames (*transformation equivariance*). Extensive experiments demonstrate that our model significantly outperforms previous methods on two sound localization benchmarks, namely, Flickr-SoundNet and VGG-Sound. Additionally, we also evaluate audio

retrieval and cross-modal retrieval tasks. In both cases, our self-supervised models demonstrate superior retrieval performances, even competitive with the supervised approach in audio retrieval. This reveals the proposed framework learns strong multi-modal representations that are beneficial to sound localisation and generalization to further applications. *The project page is https://jinxiang-liu.github.io/SSL-TIE/ .*

## CCS CONCEPTS

• **Information systems** → **Multimedia information systems**; • **Computing methodologies** → **Computer vision**; **Unsupervised learning**.

## KEYWORDS

self-supervised representation learning, sound localisation

## 1 INTRODUCTION

When looking around the world, we can effortlessly perceive the scene from multi-sensory signals, for example, whenever there is sound of dog barking, we would also expect to see a dog somewhere in the scene. A full understanding of the scene should thus include the interactions between the visual appearance and acoustic characteristics. In the recent literature, researchers have initiated research on various audio-visual tasks, including audio-visual sound separation [11–14, 46, 51–53], visual sound source localisation [6, 21–23, 30, 37, 39, 41] and audio-visual video understanding [15, 24, 28, 29, 45, 48, 50]. In this paper, we focus on the task of visual sound source localisation, with the goal to highlight the salient object by its emitted sound in a given video frame. To avoid the laborious annotations, we here consider a self-supervised setting, which only requires raw videos as the training data, *i.e.* without using any extra human annotations whatsoever.

Generally speaking, the main challenge of visual sound localisation is to learn joint embeddings for visual and audial signals. To this end, various attempts have been made in early works. [2, 39] train classification models to predict whether audio and video frame are corresponding or not. And the localisation representation is obtained by computing similarity between audio and image representations, revealing the location of sounding objects; Qian *et al.* [37] also learn audio and visual representations with the classification model to localise sounding objects, they leverage the pre-trained classifiers to aggregate more audio-image pairs of the same semantics by comparing their category labels. More recent work [6] has tried to explicitly mine the sounding regions automatically through differentiable thresholding, and then self-train the model with the InfoNCE loss [47]. Despite tremendous progress has been made, previous visual sound source localisation approaches have always neglected the important role of *data augmentations*, which has shown to be essential in self-supervised representation learning [8, 9, 16, 17].

Herein, we introduce a simple self-supervised framework to explore the efficacy of data transformation. Specifically, we exploit Siamese networks to process two different augmentations of the audio-visual pairs, and train the model with contrastive learning and geometrical consistency regularization, *i.e.* encouraging the audio-visual correspondence to be *invariant* to various transformations, while enforcing the localised sound source to be *equivariant* to geometric transformations. To validate the effectiveness of the proposed idea, we experiment with two prevalent audio-visual localisation benchmarks, namely, Flickr-SoundNet and VGG Sound-Source. Under the self-supervised setting, our approach demonstrates state-of-the-art performance, surpassing existing approaches by a large margin, even with less than 1/14 training data, thus being more data-efficient. Additionally, we also measure the quality of learned representations by two different retrieval tasks, *i.e.* audio retrieval and audio image cross-modal retrieval, which demonstrates the powerful representation learning ability of the proposed self-supervised framework.
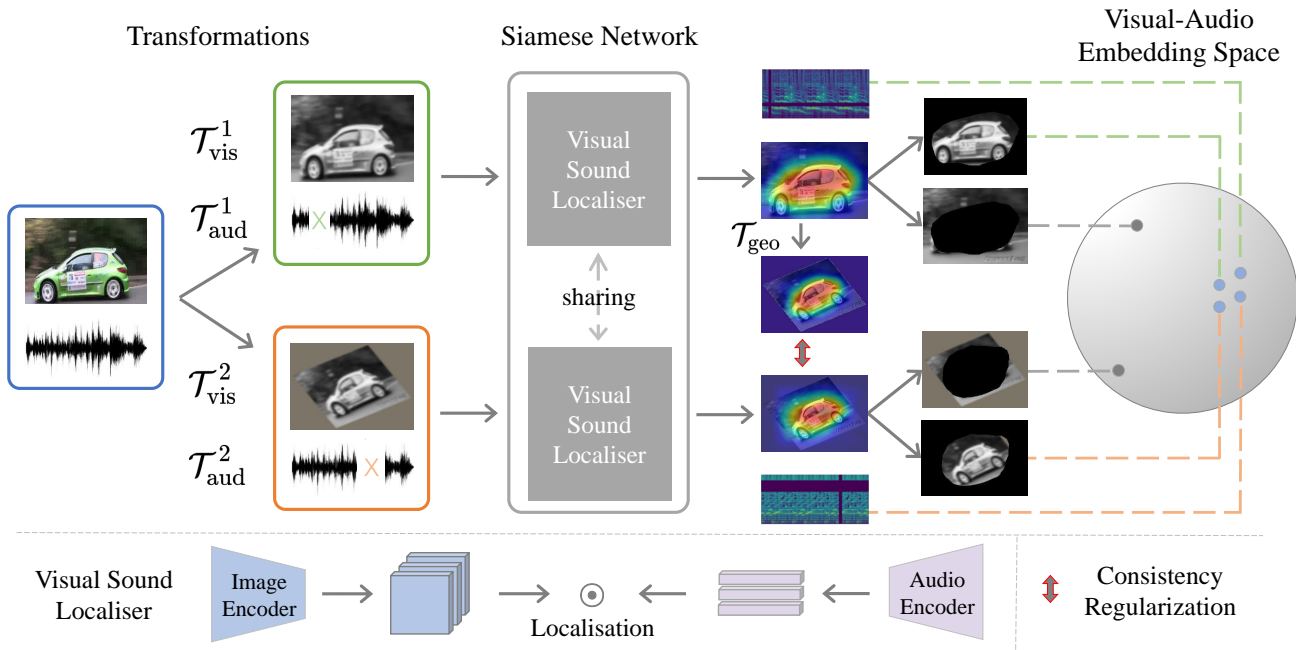
To summarise, our main contributions are three-fold: (i) We introduce a simple self-supervised framework to explore the efficacy of data transformation for visual sound localisation, concretely, we optimise a Siamese network with contrastive learning and geometrical consistency; (ii) We conduct extensive experiments and thorough ablations to validate the necessity of different augmentations, and demonstrate state-of-the-art performance on two standard sound localisation benchmarks while being more data-efficient; (iii) We initiate two audio retrieval benchmarks based on VGGSound, and demonstrate the usefulness of learned representations, *e.g.* audio retrieval and cross-modal retrieval. In both cases, our method shows impressive retrieval performances. Codes and dataset splits will be publicly released to facilitate future research.

## 2 RELATED WORK

In this section, we first review previous work on audio-visual sound source localisation, especially on the self-supervised methods; we then describe the research on self-supervised representation learning with Siamese networks; finally, we summarize the literature regarding transformation equivariance.

### 2.1 Self-supervised Sound Localisation

Audio-visual sound source localisation aims to localise the object region that corresponds to the acoustic sound in a given video frame. Early approaches have exploited the statistical models to maximize the mutual information between different modalities [10, 18]. Recently, deep neural networks have been adopted for representation learning, by leveraging the innate synchronization between audio and video, for example, SSMF [34] and AVTS [26] deploy networks to predict whether visual contents and audio are temporally aligned or not, then the sounding objects can be discovered through Class Activation Mapping (CAM) [54]. Senocak *et al.* [38] develop a foreground attention mechanism with the triplet loss [19], where the attention map is computed by the inner dot product between sound and visual context. Qian *et al.* [37] propose a two-stage framework for multiple-object sound localization, they first leverage the pre-trained classifiers to obtain pseudo category labels and then align

**Figure 2: Framework Overview. We exploit a Siamese network, with two identical branches, each branch consists of an image encoder and an audio encoder. For the one branch, we perform transformations $\mathcal{T}_{\mathbf{vis}}^1 + \mathcal{T}_{\mathbf{aud}}^1$, while for the other branch, we use transformations $\mathcal{T}_{\mathbf{vis}}^2 + \mathcal{T}_{\mathbf{aud}}^2$. In this figure, $\mathcal{T}_{\mathbf{vis}}^1$ only includes appearance transformation $\mathcal{T}_{\mathbf{app}}$, while $\mathcal{T}_{\mathbf{vis}}^2$ includes both appearance and geometric transformations $\mathcal{T}_{\mathbf{app}} + \mathcal{T}_{\mathbf{geo}}$. Both audio transformations are $\mathcal{T}_{aud}$. The framework is optimised by encouraging the audio-visual representation to be invariant to $\mathcal{T}_{app}$ and $\mathcal{T}_{geo}$, while being equivalent to $\mathcal{T}_{geo}$.**

the multi-modal features. Such pipeline is not end-to-end trainable, thus may hinder the performance.

Recently, contrastive learning with infoNCE loss [47] has shown great success in self-supervised representation learning [8, 17]. The methods including SimCLR [8] and MoCo [17] construct various augmentations of the same samples as positive pairs, while the augmentations of other samples as the negatives, resembling an instance discrimination task. Inspired by this, Chen *et al.* [6] introduce the infoNCE contrastive learning to sound source localisation, where they treat the responses of the sounding object within the foreground image with its corresponding audio as positive, while the responses of background image with audio and the responses of mismatched image-audio pairs as negatives. However, the authors ignore the importance of image data augmentations, which have proven to be critical in the self-supervised instance discrimination models [8, 9, 16, 17]. In this paper, we intend to fill this gap by exploring various data transformations.

### 2.2 Siamese Network

The Siamese network, which consists of two or more identical sub-networks, is typically used to compare the similarity between predictions brought by different entities. It is prevalent to solve many problems, including face verification [42], visual tracking [4, 27], one-shot object recognition [25], and recommendation [31]. More recently, the Siamese network has been widely adopted for

self-supervised representation representation learning [8, 9, 16, 17]. Concretely, the contrastive learning methods, such as SimCLR [8] and MoCo [17], aim to attract two augmented views of the same image while push away views from different image samples with the InfoNCE loss, thus resembling an instance discrimination loss. BYOL [16], SimSiam [9] and ContrastiveCrop [36] feed two branches of Siamese networks with different augmentations of the same image sample, and they utilize one branch to predict the output of the other. To the best of our knowledge, this is the first exploration to leverage the Siamese networks for sound localisation based on the contrastive learning.

### 2.3 Equivariant Transformation

Equivariant transformation refers that the predictions from a model are equivariant to the transformations applied to the input images. It is a popular technique in many problem which requires spatial prediction such as unsupervised landmark localisation [43, 44]. The assumption [43, 44] is that the learned landmark should be consistent with the visual effects of image deformations such as viewpoint change or object deformation. The transformation equivariance is also prevalent for some problems in semi-supervised settings including landmark localisation [20, 32], image segmentation [49], image-to-image translation [33]. The common approach of [20, 32, 33, 49] is to train the models with the labelled data and enforce the predictions for the unlabelled data to be equivariant

to the transformations applied on them. In this paper, we exploit the transformation equivariance property by integrating it into the proposed unified self-supervised framework for sound localization.

## 3 METHOD

In this paper, we consider the self-supervised audio-visual representation learning, to localise the sound source in the video frames. In Section 3.1, we first introduce the general problem scenario; In Section 3.2, we introduce the proposed Siamese framework (Figure 2), and describe different data transformations for both audio and visual signals; Lastly, in Section 3.3, we propose the essential transformation invariance and equivariance, and also summarize the training objectives for joint model optimisation.

### 3.1 Problem Scenario

In visual sound localisation, we are given a set of raw videos $\mathcal{X} = \{(I_1, A_1), (I_2, A_2), \cdots, (I_N, A_N)\}$, where $I_i \in \mathbb{R}^{3 \times H_v \times W_v}$ refers to the central frame of $i$-th video, $A_i \in \mathbb{R}^{1 \times H_a \times W_a}$ denotes its corresponding audio spectrogram, $H_v, W_v$ and $H_a, W_a$ are the spatial resolutions of two modalities respectively. The goal is to learn a visual localisation network that takes the audio-visual pair as inputs and outputs the localisation map for sounding object:

$$\Phi_{\text{loc}}(I_i, A_i; \Theta) = \mathbf{M}_{\text{loc}} \in \{0, 1\}^{H_v \times W_v} \tag{1}$$

where $\Theta$ represents the learnable parameters, and $\mathbf{M}_{\text{loc}}$ refers to a binary segmentation mask, with 1 denoting the visual location of objects that emit the sound.

### 3.2 Visual Sound Localisation

In order to learn the joint audio-visual embedding, we here exploit a Siamese network with two identical branches. As shown in Figure 2, each branch is consisted of an image encoder ($f_v(,: \theta_v)$) and an audio encoder ($f_a(,: \theta_a)$), and the embeddings of two modalities can be computed as follows:

$$\begin{aligned} v &= f_v(\mathcal{T}_{\text{vis}}(I), \theta_v), \quad v \in \mathbb{R}^{c \times h \times w} \\ a &= f_a(\mathcal{T}_{\text{aud}}(A), \theta_a), \quad a \in \mathbb{R}^c, \end{aligned} \tag{2}$$

where $\mathcal{T}_{\text{vis}}$ and $\mathcal{T}_{\text{aud}}$ refer to the augmentations imposed on visual frames and audio spectrograms, respectively. $h, w$ refer to the visual spatial resolution of the visual feature map, and $c$ denotes the dimension of the encoded audio vector.

To localise the visual objects, we can thus compute the response map $S_{i \to j}$, by measuring the cosine distance between the audio features $a_i$ and pixel-level visual features $v_j$:

$$S_{i \to j} = \frac{\langle a_i, v_j \rangle}{\|a_i\| \cdot \|v_j\|} \in \mathbb{R}^{h \times w}, \tag{3}$$

where $S_{i \to j}$ indicates the visual-audio activation between the $i$-th video frame and the $j$-th audio. The final segmentation map $\mathbf{M}_{\text{loc}}$ is attained by simply thresholding $S_{i \to j}$.

#### 3.2.1 Transformation on audio spectrogram ($\mathcal{T}_{\text{aud}}$).
Here, before feeding audio data to the audio encoder, we pre-process the 1-D waveform to obtain 2-D mel-spectrograms, with horizontal and vertical axes representing time and frequency, respectively. Then, we consider two different types of audio augmentations, i.e. spectrogram masking $\mathcal{T}_{\text{mask}}$ and audio mixing $\mathcal{T}_{\text{mix}}$.

As for spectrogram masking, we randomly replace the 2-D mel-spectrograms with zeros along two axes with random widths, that is, time masking and frequency masking on mel-spectrograms [35]. While for audio mixing, we aim to blend the audio samples with same semantics. To find the semantic identical audio for each audio sample, we compute the similarity of embedding with all other audio samples in datasets and adopt the most similar one to mix. We conduct such mixing strategy in a curriculum learning manner: the blending weights for the sampled audios are linearly increased from 0 to 0.65 as the training proceeds. Mathematically:

$$A_i^{\text{mix}} = (1 - \alpha) \cdot A_i + \alpha \cdot A_i^{\text{sim}}, \tag{4}$$

where $A_i^{\text{sim}}$ is the most similar audio sample of the audio $A_i$, $A_i^{\text{mix}}$ refers to the mixed audio, and $\alpha$ is the mixing coefficient, which increases linearly with the training epoch. In Section 4.2.3, we have conducted thorough experiments, showing both transformations are critical for improving sound localisation performance while preventing the model from overfitting.

#### 3.2.2 Transformation on visual frames ($\mathcal{T}_{\text{vis}}$).
Here, we split the image transformations into two groups: appearance transformations $\mathcal{T}_{\text{app}}$ and geometrical transformations $\mathcal{T}_{\text{geo}}$. $\mathcal{T}_{\text{app}}$ refers to transformations that only change the frame appearances, including color jittering, gaussian blur, and grayscale; $\mathcal{T}_{\text{geo}}$ changes the geometrical shapes and locations of the sounding objects, including cropping and resizing, rotation, horizontal flipping. These transformations are shown to be essential for representation learning in recent visual self-supervised approaches, e.g. SimCLR [8], MOCO [17], DINO [5], etc. We refer the readers for both audio and visual frame transformations in supplementary materials.

### 3.3 Training Details

In this section, we describe how to exploit different data transformations for training visual sound localisation models.

#### 3.3.1 Correspondence Transformation Invariance.
Though various transformations are applied on inputs, the audio-image correspondence is not altered, which means the correspondence are invariant to the transformations. Thus we still adopt batch contrastive learning for both branches in the Siamese framework to exploit the correlation between audio-visual signals, as follows:

$$m_i = \text{sigmoid}((S_{i \to i} - \epsilon)/\tau) \tag{5}$$

$$P_i = \frac{1}{|m_i|} \langle m_i, S_{i \to i} \rangle \tag{6}$$

$$N_i = \sum_{i \neq j} \frac{1}{hw} \langle \mathbf{1}, S_{i \to j} \rangle + \frac{1}{|1 - m_i|} \langle 1 - m_i, S_{i \to i} \rangle \tag{7}$$

$$\mathcal{L}_{\text{cl}} = -\frac{1}{B} \sum_{i=1}^{B} \left[ \log \frac{\exp(P_i)}{\exp(P_i) + \exp(N_i)} \right] \tag{8}$$

Here, $m_i \in \mathbb{R}^{h \times w}$ refers to the foreground pseudo-mask; $P_i$ denotes the positive set that is constructed by the responses within the mask; $N_i$ denotes the negative set, with two components: the responses between unpaired audio-visual signals and the responses of its own background.

*3.3.2* **Geometric Transformation Equivariance.** Despite the fact that $\mathcal{T}_{geo}$ on images do not change the semantic correspondences with audios, $\mathcal{T}_{geo}$ do change the predicted localisation result. And ideally, the localisation results should take the same geometrical transformations as the input images experienced during the data transformation. Formally:

$$\Phi_{loc}(\mathcal{T}_{geo}(I), A) = \mathcal{T}_{geo}(\Phi_{loc}(I, A)), \tag{9}$$

where $\Phi_{loc}(\cdot)$ refers to the sound source localisation network, and $(I, A)$ denotes the frame-audio pair.

Based on this transformation equivariance property, we implement a geometrical transformation consistency between response outputs from two branches of the Siamese framework as:

$$\mathcal{L}_{geo} = \left\| S_{i \to i}^2 \left( \mathcal{T}_{geo}(I), A \right) - \mathcal{T}_{geo}(S_{i \to i}^1(I, A)) \right\|_2, \tag{10}$$

where $S_{i \to i}^1, S_{i \to i}^2$ are response maps from the two branches of the Siamese framework, and $\|\cdot\|$ refers to the $L^2$ norm.

*3.3.3* **Optimisation Objectives.** We train the Siamese framework by jointly optimising the contrastive loss and geometrical consistency loss in a self-supervised manner,

$$\mathcal{L}_{total} = \mathcal{L}_{cl}^1 + \mathcal{L}_{cl}^2 + \lambda_{geo} \mathcal{L}_{geo}, \tag{11}$$

where $\mathcal{L}_{cl}^1, \mathcal{L}_{cl}^2$ refer to the contrastive loss in both branches, $\lambda_{geo}$ represents the weighs of $\mathcal{L}_{geo}$ and is set to 2.0 empirically.

## 4 EXPERIMENT

In this section, we conduct extensive experiments for audio-visual sound localisation on two standard benchmarks and compare with existing state-of-the-art methods. We conduct thorough ablation studies to validate the necessity of different transformations. Additionally, based on the VGGSound dataset, we introduce two new evaluation protocols on retrievals, to further evaluate the quality of learnt audio-visual representation.

### 4.1 Implementation details

Our proposed method is implemented with PyTorch. The input images are all resized to $224 \times 224$ spatial resolution, with random augmentations, including color jitterings, *e.g.* grayscale, brightness, contrast, saturation, and geometric transformations, *e.g.* rotation, horizontal flipping. For the visual and audio encoders, we here adopt the ResNet-18 backbone. The shape of output features from the visual encoder is $14 \times 14 \times 512$, and the shape of audio features is $1 \times 512$. The model is optimised with Adam using a learning rate of $10^{-4}$, and the batch size is set to 32. We train the model for 80 epochs on single GeForce RTX 3090 GPU.

### 4.2 Visual Sound Localisation

*4.2.1* **Datasets.** We train and evaluate on two datasets,

**Flickr-SoundNet [3]** contains more than 2M unconstrained image-audio pairs. Following the convention [38], we adopt a random subset with 144k image-audio sample pairs, and a subset with 10k random samples for training, termed as **Flickr-144k** and **Flickr-10k** respectively. For evaluation, we use the 250 random sampled pairs from the subset of 5k annotated sample pairs. In the test subsets,

each audio-image pair contains an image and its corresponding 20-second-long audio, which is annotated by three different subjects for reliability. The annotation is the bounding box of the dominant object that emits the sound.

**VGGSound [7]** is a video dataset consisting of 200k videos, spanning 309 sounding categories. Similar to Flickr-SoundNet splits, we use the subsets **VGGSound-144k** and **VGGSound-10k** for training. While for evaluation, we employ the VGG-SS [6], a recently-released standard testing subset from VGG-Sound dataset. The testing dataset contains 5k videos, each video is annotated with a bounding box for sounding objects in the center frame.

*4.2.2* **Metrics.** We quantitatively measure sound source localisation performance with two metrics: (i) Consensus Intersection over Union (cIoU) [38] measures the localisation accuracy through the intersection and union between ground-truth and prediction; (ii) Area Under Curve (AUC) indicates the area under the curve of cIoU plotted by varying the threshold from 0 to 1. For both metrics, high values mean better localisation performances.

*4.2.3* **Ablation Study.** In this section, we conduct thorough ablation studies on VGG-SS test, to validate the effectiveness of each component. The results are reported in Table 1. To facilitate comparisons, model-A is set as the baseline with only contrastive loss $\mathcal{L}_{cl}$ applied, which shares similar setting as LVS [6].

**Effectiveness of aggressive augmentations.** When comparing with the baseline, model-B (only appearance augmentation) and model-C (both appearance and geometrical augmentations) have clearly shown superior performance, about 3% cIoU, demonstrating the effectiveness of visual augmentations. Additionally, when adding audio augmentations (model-D), We observe further performance boost (about 4.5% cIoU over baseline).

**Effectiveness of audio mixing.** On the one hand, comparing model-D and model-E, the proposed audio mixing also brings tiny performance boost. On the other hand, we do observe its benefits for mitigating overfitting issue, as demonstrated in Figure 6 of A.2 in the appendix. For the model without leveraging audio mixing transformations, the validation cIoU tends to decrease after 40 Epochs, which is a typical performance degradation caused by severe overfitting. For the model with the audio mixing transformation, the validation loss is constantly decreasing, showing that the overfitting issue is well solved. In conclusion, our proposed audio mixing transformation can slightly improves localisation performance, as well as preventing the model from overfitting.

**Effectiveness of geometrical consistency.** When training model-F with geometrical consistency, our best model achieves the best performance, about 6% cIoU over the baseline model.

**Summary.** As shown in Table 1, all the components including various data augmentation, *e.g.* appearance and geometrical ones on visual frames, masking, and audio mixing, are all critical to boosting performance on self-supervised sound source localisation. Additionally, by further enforcing the audio-visual representation to be equivariant, the proposed framework has achieved the best performance.

**Table 1: Ablation study on the VGG-SS test set. All the models are trained with VGGSound-144k dataset. The results shows that, all data transformations and optimization losses are essential. By encouraging audio-visual invariant to various transformations, while visually equivariant to geometric transformations, we achieve considerable performance gains.**

| Model | Transformations | | | | Objectives | | Results | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{T}_{app}$ | $\mathcal{T}_{geo}$ | $\mathcal{T}_{mask}$ | $\mathcal{T}_{mix}$ | $\mathcal{L}_{cl}$ | $\mathcal{L}_{geo}$ | cIoU | AUC |
| A | | | | | ✓ | | 0.3292 | 0.3744 |
| B | ✓ | | | | ✓ | | 0.3364 | 0.3721 |
| C | ✓ | ✓ | | | ✓ | | 0.3580 | 0.3847 |
| D | ✓ | ✓ | ✓ | | ✓ | | 0.3748 | 0.3887 |
| E | ✓ | ✓ | ✓ | ✓ | ✓ | | 0.3766 | 0.3937 |
| F | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **0.3863** | **0.3965** |

**Table 2: Comparisons on the Flcikr-SoundNet test set. All the models are trained on Flickr-144k or Flickr-10k subsets. Our method significantly outperforms these competitors.**

| Method | Training Set | cIoU | AUC |
|---|---|---|---|
| Attention [38] | Flickr-10k | 0.436 | 0.449 |
| CoarseToFine [37] | Flickr-10k | 0.522 | 0.496 |
| AVO [1] | Flickr-10k | 0.546 | 0.504 |
| LVS [6] | Flickr-10k | 0.582 | 0.525 |
| **Ours** | Flickr-10k | **0.755** | **0.588** |
| Attention [38] | Flickr-144k | 0.660 | 0.558 |
| DMC [21] | Flickr-144k | 0.671 | 0.568 |
| LVS [6] | Flickr-144k | 0.699 | 0.573 |
| HPS [40] | Flickr-144k | 0.762 | 0.597 |
| SSPL [41] | Flickr-144k | 0.759 | 0.610 |
| **Ours** | Flickr-144k | **0.815** | **0.611** |

**Table 3: Comparisons on the VGG-SS and Flickr-SoundNet test sets. Note that, all models are trained on VGG-Sound 144k. For VGG-SS, our method significantly surpasses previous state-of-the-art models; when evaluating on Flickr-SoundNet test set, our method still performs the best, revealing strong *generalisation* across different datasets.**

| Method | VGG-SS | | Flickr-SoundNet | |
|---|---|---|---|---|
| | cIoU | AUC | cIoU | AUC |
| Attention [38] | 0.185 | 0.302 | 0.660 | 0.558 |
| AVO [1] | 0.297 | 0.357 | – | – |
| SSPL [41] | 0.339 | 0.380 | 0.767 | 0.605 |
| LVS [6] | 0.344 | 0.382 | 0.719 | 0.582 |
| HPS [40] | 0.346 | 0.380 | 0.768 | 0.592 |
| **Ours** | **0.386** | **0.396** | **0.795** | **0.612** |

*4.2.4* **Compare with State-of-the-Art.** Here, we compare with the existing methods on the task of sound source localisation, including: Attention [38], AVO [1], DMC [21], HPS [40], SSPL [41], CoarseToFine [37], and LVS [6].

**Quantitative Results on Flickr-SoundNet.** In Table 2, we present the comparisons between various approaches on Flickr-SoundNet test set. Here, we train the model on two training sets, namely, Flickr-10k and Flickr-144k subsets. Experimentally, our proposed method outperforms all existing methods by a large margin. Note that, some methods use additional data or information, e.g., Attention [38] uses 2796 bounding box annotated audio-image pairs as supervision. CoarseToFine [37] exploits a pretrained object detector to obtain pseudo category labels. In contrast, our proposed model is trained from scratch. Moreover, it can be seen that our model trained on 10k subset performs even better than LVS trained on 144k subset, that is to say, we achieve superior results with less than 1/14 of training data that the counterpart method [6] requires, demonstrating the high data-efficiency of our proposed framework.

**Quantitative Results on VGG-SS.** Following [6], we here train the model on the VGGSound-144k training split, but make comparisons between various approaches on the VGG-SS and Flickr-SoundNet test sets, as shown in Table 3. On VGG-SS test set, our framework surpasses the previous state-of-the-art model [40] by a noticeable margin. In addition, when evaluating on Flickr-SoundNet test set, our method also maintains its top position, revealing strong *generalisation* across different datasets.

**Open Set Sound Localisation on VGG-SS.** Following the evaluation protocol in LVS [6], in this section, we also show the sound localisation results in an open set scenario, where models are trained with 110 heard categories in VGGSound, and then evaluated on 110 heard and 110 unheard categories separately in the test set. As shown in Table 4, both approaches have experienced performance drop on unheard categories, however, our proposed model still maintains high localisation accuracy in this open set evaluation.

*4.2.5* **Qualitative Results.** In Figure 3, we show some qualitative comparisons between LVS [6] and our proposed method on Flickr-Sound test set and VGG-SS test set. As can be observed, our model generally produces more accurate localisation results than LVS, in two aspects: 1) our predictions tend to be more complete and highly consistent with the shape of the sounding objects, that means, a more precise prediction on the object boundaries, while LVS only localises the parts of objects. 2) our localisation more focuses on the foreground sounding objects, regardless of the background or silent

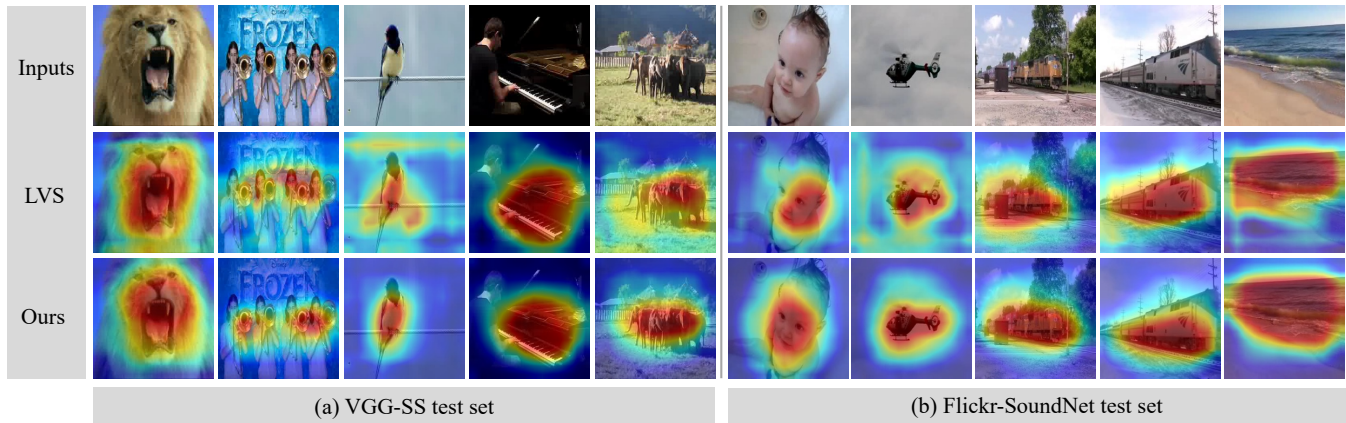(a) VGG-SS test set      (b) Flickr-SoundNet test set

**Figure 3: Qualitative results on VGG-SS and Flickr-SoundNet test sets for visual sound localisation. LVS [6], as the state-of-the-art competitor, is chosen for comparison. The models are trained on Flickr-144k and VGGSound-144k datasets respectively. Our method localises sounding objects more accurately than LVS, especially for small-size objects.**

**Table 4: Results for open set sound localisation. All models are trained on 70k samples from 110 object categories in VG-GSound, and evaluated on 110 heard categories and 110 unheard categories. Our method shows strong performance.**

| Test class | Method | CIoU | AUC |
|---|---|---|---|
| Heard 110 | LVS [6] | 0.289 | 0.362 |
| | Ours | **0.390** | **0.403** |
| Unheard 110 | LVS [6] | 0.263 | 0.347 |
| | Ours | **0.365** | **0.386** |

**Table 5: Results for audio retrieval. For fair comparisons, all models adopt the ResNet-18 backbone. We here use Accuracy (A@5, A@10) and Precision (P@1, P@5) as metrics. Our learned audio representations are powerful and sometimes comparable to full supervision.**

| Method | Supervision | A@5 | A@10 | P@1 | P@5 |
|---|---|---|---|---|---|
| Random | No | 20.10 | 28.06 | 13.88 | 6.06 |
| VGG-H [7] | Full | **42.07** | **45.27** | 58.69 | **27.63** |
| LVS [6] | Self | 26.01 | 33.67 | 21.17 | 9.37 |
| **Ours** | Self | 41.15 | 44.19 | **60.19** | 27.55 |

distracting objects; while the localisations of LVS are sometimes scattered even in the clean background, *e.g.*, the *third* column in subplot (a) and the *second* and *third* column in subplot (b).

## 4.3 Audio Retrieval

To further investigate the quality our our learned audio representation, we evaluate the methods on audio retrieval task.

*4.3.1* **Benchmarks.** Due to the lack of unified benchmarks, we first divide the VGGSound dataset into train-val set and test set, with categories being disjoint. The former is for training and validation,

**Table 6: Results for audio-image cross-modal retrieval. We report Accuracy (A@5, A@10) and Precision (P@1, P@5). Our model has shown impressive retrieval performance, implying the strong multi-modal representation extraction abilities of our self-supervised models.**
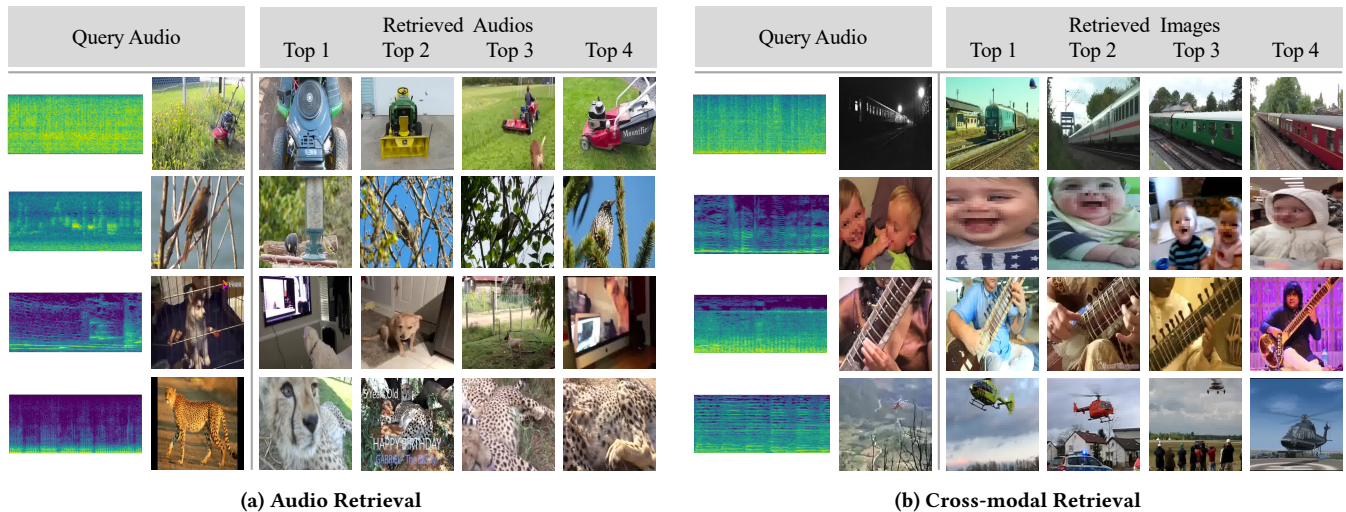
| Method | Train Category | A@5 | A@10 | P@1 | P@5 |
|---|---|---|---|---|---|
| Random | 0 | 4.44 | 8.01 | 1.35 | 1.61 |
| LVS [6] | All | 11.00 | 16.26 | 10.48 | 4.93 |
| Ours | 110 | 22.44 | 27.72 | 25.50 | 12.34 |
| Ours | All | **31.67** | **35.81** | **40.91** | **19.52** |

while the latter consisting of unseen categories is for evaluation. In detail, the train-val set spans over 274 categories with 169923 samples, we randomly sample a 144k subset for training and the rest as the validation set. The test set has 35 object categories covering 20304 samples.

*4.3.2* **Metrics.** We use two standard metrics: *accuracy* and *precision*. For Top-K accuracy (A@K), as long as the K results contain at least one item of the same category as the query audio, the retrieval is regarded as correct. Precision (P@K) is the percentage of the top-K retrieved items of the same category with query audio.

*4.3.3* **Baselines.** Here, we compare the retrieval results with the following models: 1) Random: the model weights are randomly initialized without training. 2) VGG-H: the model is trained with ground-truth category supervision on the training set, as has been done in [7], 3) LVS: a recent state-of-the-art model trained for visual sound localisation [6]. 4) Ours: our Siamese framework trained on self-supervised visual sound localisation. For fair comparisons, all models use the ResNet-18 backbone as the audio encoder.

*4.3.4* **Retrieval Detail.** For each query audio in the test set, we extract the *512-D* feature with the audio encoder from different models, *e.g.* baselines and our model; we then calculate the cosine

**(a) Audio Retrieval**



**(b) Cross-modal Retrieval**

**Figure 4: Qualitative results on retrieval tasks. (a) Audio retrieval, which retrieves semantic-identical audios with the query audio. (b) Audio image cross-modal retrieval, which we use the audio as query to retrieve images. The results show our model can accurately retrieve samples with close semantics, indicating our framework has learnt powerful multi-modal representation.**

similarity between the query audio and all the rest samples; finally, we rank the similarity in a descending order, and output the top-K retrieved audios.

*4.3.5* **Results.** We report the results in Table 5, as can be observed, our self-supervised model significantly outperforms the random and LVS baselines and even demonstrate comparable results to the fully-supervised model, *i.e.* (VGG-H). In Figure 4 (a), we qualitatively show some audio retrieval results in the form of paired video frames. Our model can correctly retrieve samples with close semantics, which can potentially be used as auxiliary evidence for video retrieval applications.

## 4.4 Cross-modal Retrieval

We also evaluate an audio-image cross-modal retrieval task to evaluate the learned cross-modal representations.

*4.4.1* **Benchmark.** Similar to Section 4.3.1, we obtain the train set and test set from VGGSound dataset. The test set has 20304 samples spanning 35 categories which are the same as audio retrieval. The train sets have two versions which both have 144k samples. The difference is one train set covers all categories while the other train set has 110 categories which are disjoint with test set.

*4.4.2* **Metrics.** Similar to the audio retrieval task, we also report Top-K accuracy (A@K) and Top-K precision (P@K).

*4.4.3* **Baselines.** We compare the retrieval results with the following models: 1) Random 2) LVS [6] 3) Ours. For fair comparisons, all models employ the ResNet-18 backbone as audio and image encoders. See Section 4.3.3 for more details.

*4.4.4* **Retrieval Details.** For each query audio in the test set, we extract *512-D* feature with the audio encoder from different models. For all images to be retrieved in the dataset, we extract the visual features from the visual encoder and spatially pool them into *512-D*

vector. Then we compute the cosine similarity between the query audio and the image samples to be retrieved. Finally, we rank the similarity in descending order and check the category labels from top-K retrieved images.

*4.4.5* **Results.** We report the cross-modal retrieval results in Table 6. Comparing with baselines, our representations from self-supervised sound localiser achieve impressive cross-modality retrieval performances, without any finetuning. We also qualitatively show the results in Figure 4 (b). The quantitative and qualitative results show that the various transformations in the proposed sound localisation framework have enabled the audio and visual encoders very strong representation abilities. As a result, our self-supervised framework is remarkably effective for sound source localisation as well as multi-modal retrieval tasks.

## 5 CONCLUSION

This paper has presented a self-supervised framework for sound source localisation, by fully exploiting various transformations. The motivation is that appearance transformations and geometrical transformations on image-audio pairs are coming with two implicit but significant properties: *invariance* and *equivariance*. Invariance refers that the audio-image correspondences are invariant to data transformations; while equivariance denotes the localisation results are equivariant to the geometrical transformations that applied to input images. Combining these, we propose Siamese networks with dual branches, each branch accepts input data with different transformations on both modalities. Thanks to the two properties, the framework is trained in a fully self-supervised way. Experiments demonstrate our method significantly outperforms current methods in visual sound localisation. Additionally, we also evaluate audio retrieval and cross-modal retrieval tasks, to show the learned powerful multi-modal representations. Finally, we perform thorough

ablation studies to verify the effectiveness of each component in the framework.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. 2020. Self-supervised learning of audio-visual objects from video. In *European Conference on Computer Vision*. Springer, 208–224.

[2] Relja Arandjelovic and Andrew Zisserman. 2018. Objects that sound. In *Proceedings of the European conference on computer vision (ECCV)*. 435–451.

[3] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. 2016. Soundnet: Learning sound representations from unlabeled video. In *Advances in Neural Information Processing Systems*.

[4] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. 2016. Fully-convolutional siamese networks for object tracking. In *European conference on computer vision*. Springer, 850–865.

[5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9650–9660.

[6] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. 2021. Localizing visual sounds the hard way. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16867–16876.

[7] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. 2020. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 721–725.

[8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.

[9] Xinlei Chen and Kaiming He. 2021. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15750–15758.

[10] John W Fisher III, Trevor Darrell, William Freeman, and Paul Viola. 2000. Learning joint statistical models for audio-visual fusion and segregation. *Advances in neural information processing systems* 13 (2000).

[11] Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. 2020. Music gesture for visual sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10478–10487.

[12] Ruohan Gao, Rogerio Feris, and Kristen Grauman. 2018. Learning to separate object sounds by watching unlabeled video. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 35–53.

[13] Ruohan Gao and Kristen Grauman. 2019. Co-separating sounds of visual objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3879–3888.

[14] Ruohan Gao and Kristen Grauman. 2021. Visualvoice: Audio-visual speech separation with cross-modal consistency. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 15490–15500.

[15] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. 2020. Listen to look: Action recognition by previewing audio. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10457–10467.

[16] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. 2020. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems* 33 (2020), 21271–21284.

[17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9729–9738.

[18] John Hershey and Javier Movellan. 1999. Audio vision: Using audio-visual synchrony to locate sounds. *Advances in neural information processing systems* 12 (1999).

[19] Elad Hoffer and Nir Ailon. 2015. Deep metric learning using triplet network. In *International workshop on similarity-based pattern recognition*. Springer, 84–92.

[20] Sina Honari, Pavlo Molchanov, Stephen Tyree, Pascal Vincent, Christopher Pal, and Jan Kautz. 2018. Improving landmark localization with semi-supervised learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1546–1555.

[21] Di Hu, Feiping Nie, and Xuelong Li. 2019. Deep multimodal clustering for unsupervised audiovisual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9248–9257.

[22] Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou. 2020. Discriminative sounding objects localization via self-supervised audiovisual matching. *Advances in Neural Information Processing Systems* 33 (2020), 10077–10087.

[23] Di Hu, Yake Wei, Rui Qian, Weiyao Lin, Ruihua Song, and Ji-Rong Wen. 2021. Class-aware Sounding Objects Localization via Audiovisual Correspondence. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).

[24] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. 2019. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5492–5501.

[25] Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, Vol. 2. Lille, 0.

[26] Bruno Korbar, Du Tran, and Lorenzo Torresani. 2018. Cooperative learning of audio and video models from self-supervised synchronization. *Advances in Neural Information Processing Systems* 31 (2018).

[27] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. 2016. Learning by tracking: Siamese CNN for robust target association. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 33–40.

[28] Jun-Tae Lee, Mihir Jain, Hyoungwoo Park, and Sungrack Yun. 2020. Cross-attentional audio-visual fusion for weakly-supervised action localization. In *International Conference on Learning Representations*.

[29] Yan-Bo Lin, Yu-Jhe Li, and Yu-Chiang Frank Wang. 2019. Dual-modality seq2seq network for audio-visual event localization. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2002–2006.

[30] Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, and Ming-Hsuan Yang. 2021. Unsupervised sound localization via iterative contrastive learning. *arXiv preprint arXiv:2104.00315* (2021).

[31] Saket Maheshwary and Hemant Misra. 2018. Matching resumes to jobs via deep siamese network. In *Companion Proceedings of the The Web Conference 2018*. 87–88.

[32] Olga Moskvyak, Frederic Maire, Feras Dayoub, and Mahsa Baktashmotlagh. 2021. Semi-supervised Keypoint Localization. In *International Conference on Learning Representations*. https://openreview.net/forum?id=yFJ67zTeI2

[33] Aamir Mustafa and Rafał K Mantiuk. 2020. Transformation consistency regularization–a semi-supervised paradigm for image-to-image translation. In *European Conference on Computer Vision*. Springer, 599–615.

[34] Andrew Owens and Alexei A Efros. 2018. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 631–648.

[35] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. 2019. Specaugment: A simple data augmentation method for automatic speech recognition. (2019), 2613–2617.

[36] Xiangyu Peng, Kai Wang, Zheng Zhu, and Yang You. 2022. Crafting Better Contrastive Views for Siamese Representation Learning. *arXiv preprint arXiv:2202.03278* (2022).

[37] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. 2020. Multiple sound sources localization from coarse to fine. In *European Conference on Computer Vision*. Springer, 292–308.

[38] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. 2018. Learning to localize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4358–4366.

[39] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. 2019. Learning to localize sound sources in visual scenes: Analysis and applications. *IEEE transactions on pattern analysis and machine intelligence* 43, 5 (2019), 1605–1619.

[40] Arda Senocak, Hyeonggon Ryu, Junsik Kim, and In So Kweon. 2022. Learning Sound Localization Better From Semantically Similar Samples. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE.

[41] Zengjie Song, Yuxi Wang, Junsong Fan, Tieniu Tan, and Zhaoxiang Zhang. 2022. Self-Supervised Predictive Learning: A Negative-Free Method for Sound Source Localization in Visual Scenes. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.

[42] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, and Lior Wolf. 2014. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1701–1708.

[43] James Thewlis, Samuel Albanie, Hakan Bilen, and Andrea Vedaldi. 2019. Unsupervised learning of landmarks by descriptor vector exchange. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 6361–6371.

[44] James Thewlis, Hakan Bilen, and Andrea Vedaldi. 2017. Unsupervised learning of object landmarks by factorized spatial embeddings. In *Proceedings of the IEEE international conference on computer vision*. 5916–5925.

[45] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. 2020. Unified multisensory perception: Weakly-supervised audio-visual video parsing. In *European Conference on Computer Vision*. Springer, 436–454.

[46] Efthymios Tzinis, Scott Wisdom, Aren Jansen, Shawn Hershey, Tal Remez, Dan Ellis, and John R. Hershey. 2021. Into the Wild with AudioScope: Unsupervised Audio-Visual Separation of On-Screen Sounds. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. https://openreview.net/forum?id=MDsQkFP1Aw

[47] Aaron Van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning with Contrastive Predictive Coding. *CoRR* abs/1807.03748 (2018). arXiv:1807.03748 http://arxiv.org/abs/1807.03748

[48] Weiyao Wang, Du Tran, and Matt Feiszli. 2020. What makes training multimodal classification networks hard?. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12695–12705.

[49] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. 2020. Selfsupervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12275–12284.

[50] Fanyi Xiao, Yong Jae Lee, Kristen Grauman, Jitendra Malik, and Christoph Feichtenhofer. 2020. Audiovisual slowfast networks for video recognition. *arXiv preprint arXiv:2001.08740* (2020).

[51] Xudong Xu, Bo Dai, and Dahua Lin. 2019. Recursive visual sound separation using minus-plus net. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 882–891.

[52] Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. 2019. The sound of motions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1735–1744.

[53] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. 2018. The sound of pixels. In *Proceedings of the European conference on computer vision (ECCV)*. 570–586.

[54] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2921–2929.

# A   APPENDIX

## A.1   Data Transformations

We deploy various data transformation on the input data including audios and video frames. These transformations are visualized in Figure 5.
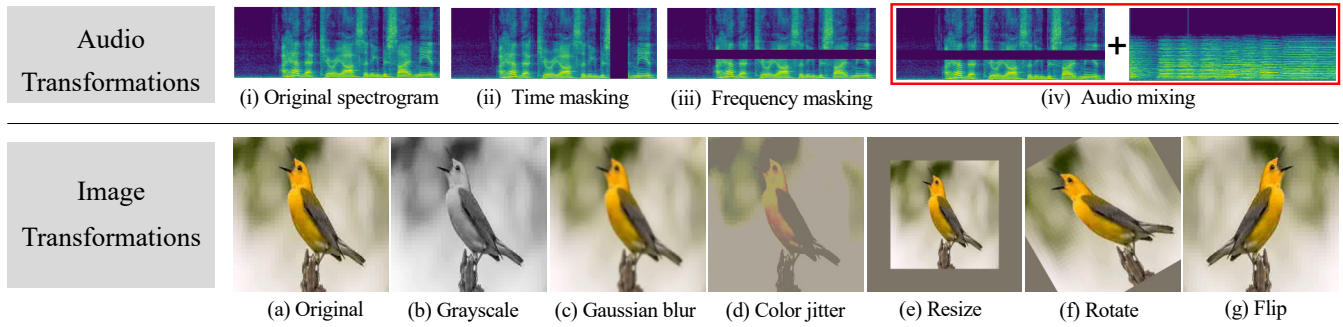
In the training stage, the audios are mixed with semantic-identical audio samples and then randomly masked with two strategies, namely time masking and frequency masking. The masking probabilities of two masking strategies are both 0.8.

For the transformations on video frames $\mathcal{T}_{vis}$: the color jitter randomly changes the brightness, contrast, saturation and hue of the image. And the strength for changing the four factors is the tuple (0.4, 0.4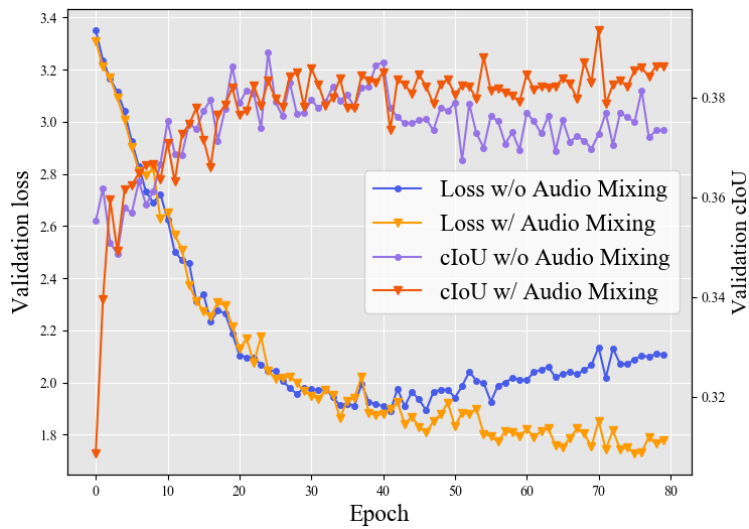, 0.4, 0.1), each element corresponding to one factor respectively; the application probability is 0.8. The grayscale transformation is applied with the probability of 0.2. For the Gaussian blur, the standard deviation for creating the blurring kernel is uniformly sample from [0.1, 2.0]; the application probability is 0.5. For the geometrical transformations $\mathcal{T}_{geo}$: the resize factor is 0.5 with the probability of 0.5. The max rotation degrees is 30. The horizontal flip operation is deployed with the probability of 0.5.

## A.2   Audio Mixing Transformation

The validation curves of VGGSound-144k with or without audio mixing are shown in Fig 6. As illustrated, the audio mixing transformation can bring performance gains by preventing the model for overfitting.

| Audio Transformations | (i) Original spectrogram | (ii) Time masking | (iii) Frequency masking | (iv) Audio mixing |

| Image Transformations | (a) Original | (b) Grayscale | (c) Gaussian blur | (d) Color jitter | (e) Resize | (f) Rotate | (g) Flip |

**Figure 5: Transformations adopted in the proposed framework. For the visual domain, we explore two types of image transformations: appearance transformation $\mathcal{T}_{app}$ (b-d) and geometrical transformation $\mathcal{T}_{app}$ (e-g); for the audial domain, we apply three effective transformations (ii-iv), which are denoted as $\mathcal{T}_{aud}$.**



**Figure 6: Validation curves with or without the audio mixing transformation on VGGSound-144k. Such audio mixing transformation can bring tiny performance boost, and prevent the model from overfitting.**