# Audio-visual Generalised Zero-shot Learning
# with Cross-modal Attention and Language

Otniel-Bogdan Mercea[1], Lukas Riesch[1,2], A. Sophia Koepke[1], Zeynep Akata[1,3,4]

[1]University of Tübingen  [2]Robert Bosch GmbH  [3]MPI for Informatics
[4]MPI for Intelligent Systems

{otniel-bogdan.mercea, a-sophia.koepke, zeynep.akata}@uni-tuebingen.de
lukas.riesch@de.bosch.com

## Abstract

*Learning to classify video data from classes not included in the training data, i.e. video-based zero-shot learning, is challenging. We conjecture that the natural alignment between the audio and visual modalities in video data provides a rich training signal for learning discriminative multi-modal representations. Focusing on the relatively underexplored task of audio-visual zero-shot learning, we propose to learn multi-modal representations from audio-visual data using cross-modal attention and exploit textual label embeddings for transferring knowledge from seen classes to unseen classes. Taking this one step further, in our generalised audio-visual zero-shot learning setting, we include all the training classes in the test-time search space which act as distractors and increase the difficulty while making the setting more realistic. Due to the lack of a unified benchmark in this domain, we introduce a (generalised) zero-shot learning benchmark on three audio-visual datasets of varying sizes and difficulty, VGGSound, UCF, and ActivityNet, ensuring that the unseen test classes do not appear in the dataset used for supervised training of the backbone deep models. Comparing multiple relevant and recent methods, we demonstrate that our proposed AVCA model achieves state-of-the-art performance on all three datasets. Code and data are available at* https://github.com/ExplainableML/AVCA-GZSL.

## 1. Introduction

Most zero-shot learning (ZSL) methods developed for image classification [6, 7, 62, 63, 75, 85] and action recognition [13, 14, 32, 87] only use unimodal input, e.g. images. However, humans leverage multi-modal sensory inputs in their everyday activities. Imagine the situation in which the sound of a dog barking is audible but the dog is visually oc-
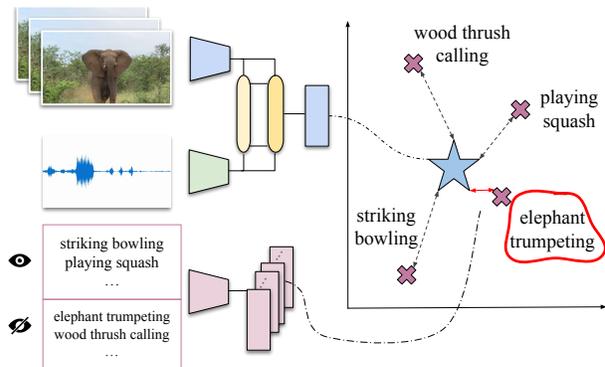


Figure 1. Our audio-visual (generalised) ZSL framework aligns an audio-visual embedding with the corresponding textual label embedding via cross-modal attention. It can classify videos from previously unseen classes (*e.g. elephant trumpeting*) by predicting the class (red) whose textual label embedding (purple cross) is closest to the audio-visual embedding (blue star).

cluded. In this case, we cannot understand the scene when relying on visual information alone. Using multiple modalities, such as vision and sound, allows to gather context and capture complementary information. Similarly, using both visual and audio information allows for a richer training signal for learning frameworks. This paper investigates the challenging task of (generalised) ZSL with multi-modal audio-visual data by leveraging the natural alignment of audio and visual information in videos.

Recently, [44, 56] have explored the task of zero-shot video recognition using multi-modal visual and audio information as inputs. However, the AudioSetZSL dataset [56] used for this, contains an overlap between the classes used for validation and testing. This results in learning stronger representations for classes overlapping with the training and validation sets (which covers all the classes in this dataset) and hinders the model's capability to learn sufficiently gen-

eralisable representations that allow information transfer. In real-world applications, such models perform well on seen classes, but poorly on previously truly unseen classes. In this work, we propose three benchmarks of varying size and difficulty curated from the VGGSound [19], UCF101 [64], and ActivityNet [26] datasets that could act as a unified and challenging playground for Generalised ZSL (GZSL) and ZSL research in the audio-visual domain. We suggest using audio and visual features extracted using SeLaVi [10] pretrained using self-supervision. Throughout this work, we use features that were obtained from training in a self-supervised fashion to reduce the information leakage from supervised pre-training to the zero-shot task which has been identified as a problem in other ZSL benchmarks [14].

We tackle the audio-visual generalised zero-shot learning task with our Audio-Visual Cross-Attention (AVCA) framework which is trained to align a rich learnt audio-visual representation with textual label embeddings. Our multi-stream architecture contains an audio and a visual branch which exchange information using cross-attention between the two modalities. AVCA is computationally lightweight and efficient since it uses audio and visual features extracted from pretrained networks as inputs instead of raw audio and image data. Our proposed framework is trained using multiple novel loss functions that are based on triplet losses and a regularisation loss that ensures that salient unimodal information is preserved in the learnt multi-modal representations. Our experiments show that AVCA achieves state-of-the-art performance on the three introduced benchmark datasets. We show that using multi-modal input data leads to stronger (G)ZSL performance than using unimodal data.

To summarise, our contributions are as follows: (1) We introduce three novel benchmarks for audio-visual (generalised) zero-shot learning curated from the VGGSound, UCF101, and ActivityNet datasets; (2) We propose AVCA, a cross-modal model for audio-visual (G)ZSL which leverages cross-modal attention between audio and visual information; (3) We show that AVCA yields state-of-the-art performance on all proposed audio-visual (G)ZSL benchmarks, outperforming the state-of-the-art unimodal and multi-modal zero-shot learning methods. Furthermore, we provide a qualitative analysis of the learnt multi-modal embedding space, demonstrating well-separated clustering for both seen and unseen classes.

## 2. Related Work

We review audio-visual learning, ZSL with image, video and audio data, and audio-visual ZSL.

**Audio-visual learning.** Audio-visual learning has enabled tremendous progress for numerous applications, such as for separating and localising sounds in videos [2, 5, 9, 18, 29, 53, 60, 69, 71, 84, 89, 90, 92], audio-visual synchronisa-

tion [17, 23, 25, 38], person-clustering in videos [15], (visual) speech and speaker recognition [3, 4, 48], spotting of spoken keywords [47, 59], audio synthesis using visual information [28,30,41,42,51,66,67,91], and audio-driven image synthesis [36, 77]. Additionally, the natural alignment between audio and visual data in videos has been leveraged to learn powerful audio-visual representations for video or audio classification [8,10,11,20,21,43,49,54,55,57,81]. In contrast to those methods, we consider the ZSL setting for classification.

**ZSL with images, videos and audio.** Recently, numerous image-based generative ZSL methods have been proposed [52, 63, 75, 79, 80, 93, 95]. Their drawback is that the unseen classes need to be known a priori. In contrast, non-generative methods [6, 7, 27, 40, 62, 78, 85, 86] learn a mapping from input features to semantics of the classes (*e.g.* textual class label embeddings). Our AVCA model also learns to map its inputs to textual embeddings, but it leverages cross-attention between the audio and visual input modalities rather than using only visual inputs.

Video-based ZSL has been addressed by multiple recent works [13, 14, 31, 32, 61, 76, 87]. Using features extracted from pretrained networks results in computationally more feasible frameworks [13, 32, 76] than training end-to-end [14]. Our model also takes pre-extracted audio and visual features as inputs, resulting in a computationally efficient framework. In order to consider a pure ZSL setting when using pre-extracted features, the overlap between classes used for supervised pre-training of the feature extractors and unseen classes has to be removed [14, 31, 61]. This was not done in some of the previous works (*e.g.* [13, 32, 76, 94]). In contrast, we propose three benchmarks for audio-visual (G)ZSL on multi-modal audio-visual video datasets with no overlap between classes used for supervised pre-training and unseen classes.

Methods for zero-shot audio classification [82, 83] also used textual sound class embeddings (*e.g.* word2vec [45], BERT [24], or GloVe [58]) or descriptions. [22] investigate zero-shot music classification and tagging with word2vec embeddings and human-labeled attribute information (e.g. the presence or absence of musical instruments). For our AVCA model, we do not use any attribute information, but instead leverage the semantic alignment between audio and visual information in addition to textual label embeddings.

**Audio-visual ZSL.** Recently, [44,56] proposed frameworks that consider the task of GZSL from audio-visual data. AVGZSLNet [44] uses late fusion on the AudioSetZSL dataset [56] to combine information from the two modalities. Instead, and also different to other audio-visual frameworks [68, 88] that use a simple dot-product operation for cross-attention, we use a transformer-based cross-attention mechanism. This allows for early and efficient sharing of multi-modal information, which is further encouraged by
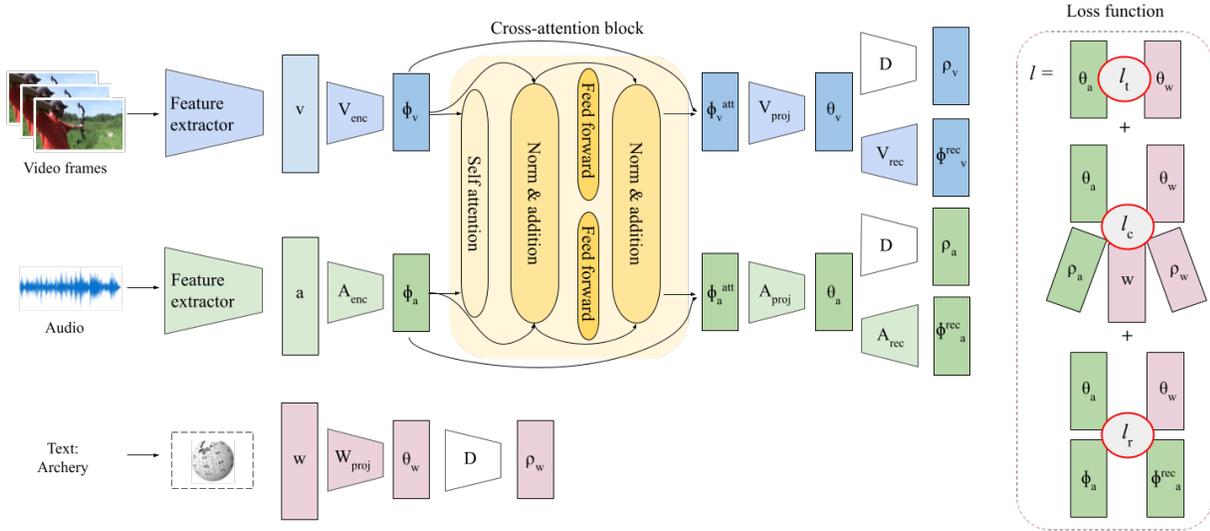
Figure 2. Our Audio-Visual Cross Attention (AVCA) model takes visual and audio features as inputs. A cross-attention block allows the sharing of information across modalities. The outputs of the two model branches are trained to be aligned with their corresponding textual label embedding using losses illustrated on the right-hand side. Negative samples for the contrastive loss functions are obtained using visual and audio inputs from different videos which do not share semantic information. We only show losses that involve the audio branch, those for the visual branch are similar. At test time, the class prediction is obtained by determining the class for which $\theta_w$ is closest to $\theta_v$.

our proposed loss functions. Furthermore, the AudioSet-ZSL dataset [56] does not include a validation split with unseen validation classes. Hence, [44, 56] select the GZSL hyperparameters directly on the (unseen) test classes. Furthermore, the AudioSetZSL dataset is comparatively small; it uses only 10 test classes as unseen classes. To allow for evaluation of audio-visual ZSL at larger scale and in a pure GZSL setting, we propose new benchmarks on three different audio-visual video datasets. Our proposed benchmarks are suitable for both the GZSL and ZSL tasks.

## 3. Audio-Visual Cross Attention (AVCA)

The goal of audio-visual ZSL from video data is to learn to recognise videos from unseen classes (U), *i.e.* classes that were not seen during training. In the GZSL setting, the test set contains not only samples from unseen classes, but also from seen classes (S). This makes GZSL more challenging and more closely aligned with real-world learning tasks.

More formally, we denote the training set consisting only of samples from seen classes by $S = (v_i^s, a_i^s, y_i^s)_{i \in \{1, \cdots, N\}}$, where $v_i^s, a_i^s$ are visual and audio features respectively, $y_i^s$ is the corresponding ground-truth class $j$, and $N$ is the number of samples in the training set. We refer to the class-level text embedding for class $j$ as $w_j^s$. The goal is to learn a function $h : (v_i^s, a_i^s) \mapsto w_j^s$ which can then also be applied to samples from unseen classes $h(v_i^u, a_i^u) = w_j^u$, where $(v_j^u, a_j^u, y_j^u) \in U$ for the set of test samples from unseen classes $U = (v_i^u, a_i^u, y_i^u)_{i \in \{1, \cdots, M\}}$ with $M$ samples.

### 3.1. Model Architecture

Our AVCA model architecture is visualised in Fig. 2. For easier readability, we dropped the subscripts $i$, $j$, indicating the $i$-th dataset sample and the ground-truth class $j$.

AVCA takes audio and visual features $a, v \in \mathbb{R}^{k_{input}}$ as inputs which are extracted using pretrained feature extractors. Those are passed through two different encoder blocks $A_{enc}$ and $V_{enc}$ for the audio and visual modality respectively, giving embeddings

$$A_{enc}(a) = \phi_a \text{ and } V_{enc}(v) = \phi_v \tag{1}$$

with $\phi_a, \phi_v \in \mathbb{R}^{k_f}$. The encoder blocks each consist of a sequence of two linear layers $f_1^m, f_2^m$ for $m \in \{a, v\}$, where $f_1^m : \mathbb{R}^{k_{input}} \to \mathbb{R}^{k_{fhidd}}$ and $f_2^m : \mathbb{R}^{k_{fhidd}} \to \mathbb{R}^{k_f}$. $f_1^m, f_2^m$ are each followed by batch normalisation [35], a ReLU [50], and dropout [65] with dropout rate $r_{enc}$.

**Cross-attention block.** We propose to use a cross-attention block to share information between the audio and visual representations $\phi_a$ and $\phi_v$. It consists of a multi-head self-attention layer, followed by a fully-connected feed-forward block. Similar to [74], we use a residual connection for the two layers, followed by layer normalisation [12].

The feed-forward blocks for the audio and visual branch each consist of a linear projection layer $f_3^m : \mathbb{R}^{k_f} \to \mathbb{R}^{k_{attnhidd}}$ for $m \in \{a, v\}$, followed by GELU [33], dropout with dropout rate of $r_{enc}$, another linear projection layer $f_4^m : \mathbb{R}^{k_{attnhidd}} \to \mathbb{R}^{k_f}$ for $m \in \{a, v\}$ and finally a dropout with dropout rate of $r_{enc}$. The outputs of the cross-attention block are $\phi_a^{att}, \phi_v^{att} \in \mathbb{R}^{k_f}$.

A residual connection around the cross-attention block and subsequent projection blocks $A_{proj}$ and $V_{proj}$ give

$$A_{proj}(\phi_a^{att} + \phi_a) = \theta_a \text{ and } V_{proj}(\phi_v^{att} + \phi_v) = \theta_v, \quad (2)$$

where $\theta_a, \theta_v \in \mathbb{R}^{k_{proj}}$. The projection blocks each consist of a sequence of two linear layers $f_5^m$ and $f_6^m$ for $m \in \{a, v\}$, where $f_5^m : \mathbb{R}^{k_f} \to \mathbb{R}^{k_{fhidd}}$ and $f_6^m : \mathbb{R}^{k_{fhidd}} \to \mathbb{R}^{k_{proj}}$. $f_5^m, f_6^m$ are each followed by batch normalisation, a ReLU, and dropout with dropout rate $r_{proj}$.

Furthermore, the word2vec class label embeddings $w^j$ for class $j$ are passed through the projection block $W_{proj}(w^j) = \theta_w^j$, where $\theta_w^j \in \mathbb{R}^{k_{proj}}$ (in Fig. 2 shown without the superscript $j$). $W_{proj}$ consists of a sequence of one linear projection layer, batch normalisation, ReLU, and dropout with dropout rate $r_{dec}$.

At test time, class predictions $c$ are obtained by determining the class $c$ that corresponds to the textual class label embedding that is closest to the multi-modal representation $\theta_v$ (in our experiments we found that using $\theta_a$ gave slightly weaker results):

$$c = \underset{j}{\arg\min}(\|\theta_w^j - \theta_v\|_2). \quad (3)$$

### 3.2. Loss Functions

We train our AVCA model using a loss function $l$ consisting of a base triplet loss $l_t$, a composite triplet and reconstruction loss $l_c$, and a regularisation loss $l_r$:

$$l = l_t + l_c + l_r. \quad (4)$$

We use the triplet loss function $t(a, p, n) = \max(\|a - p\|_2 - \|a - n\|_2 + \mu)$, where $a$ is the anchor embedding, $p$ and $n$ are embeddings for positive samples and negative samples respectively, and $\mu$ is the margin hyperparameter. For triplet losses, we use the superscript $+$ to denote positive samples that match the anchor and $-$ for negative samples that do not semantically match the anchor. For all other losses, we only use matching pairs.

**Base triplet loss.** In our base triplet loss $l_t$:

$$l_t = t(\theta_a^+, \theta_w^+, \theta_a^-) + t(\theta_v^+, \theta_w^+, \theta_v^-) \\ + t(\theta_w^+, \theta_a^+, \theta_w^-) + t(\theta_w^+, \theta_v^+, \theta_w^-), \quad (5)$$

where $\theta_m^+$ and $\theta_m^-$ correspond to positive and negative samples respectively for $m \in \{a, v, w\}$, ensuring that the projected visual and audio features $\theta_v$ and $\theta_a$ are aligned with the projected textual features $\theta_w$. This is essential, since at test time, the proximity of $\theta_v$ (which, despite being the output of the visual branch of AVCA, is a multi-modal embedding containing both audio and visual information) to $\theta_w$ for different classes is used to determine the output class.

**Composite triplet and reconstruction loss.** Inspired by [44], we additionally use a composite triplet and reconstruction loss and explain its components in more detail below:

$$l_c = l_{rec} + l_{ct} + l_w. \quad (6)$$

We use a decoder $D : \mathbb{R}^{k_{proj}} \mapsto \mathbb{R}^{k_{w2v}}$, such that $D(\theta_m) = \rho_m$ for $m \in \{a, v, w\}$. $D$ consists of a sequence of one linear projection layer, batch normalisation, a ReLU, and dropout with dropout rate $r_{dec}$. We employ the mean squared error metric $d(b, c) = \frac{1}{n}\sum_{i=1}^{n}(b_i - c_i)^2$. The reconstruction loss $l_{rec}$ can then be written as:

$$l_{rec} = d(\rho_a, w) + d(\rho_v, w) + d(\rho_w, w). \quad (7)$$

This ensures that AVCA is able to decode the pre-extracted textual label embeddings $w$ from the embeddings $\theta_a, \theta_v, \theta_w$. The triplet loss $l_{ct}$ is defined as follows:

$$l_{ct} = t(\rho_w^+, \rho_a^+, \rho_a^-) + t(\rho_w^+, \rho_v^+, \rho_v^-), \quad (8)$$

where $\rho^+$ and $\rho^-$ correspond to positive and negative examples respectively. $l_{ct}$ further encourages the decoded audio and visual features $\rho_a, \rho_v$ to be aligned with the textual features $\rho_w$ that were obtained using the same decoder (with shared weights). The third component $l_w$ of $l_c$ is similar to the base triplet loss in Eq. (5) and compares the audio and visual embeddings $\theta_a, \theta_v$ to $\theta_w$:

$$l_w = t(\theta_w^+, \theta_a^+, \theta_a^-) + t(\theta_w^+, \theta_v^+, \theta_v^-) \\ t(\theta_a^+, \theta_w^+, \theta_w^-) + t(\theta_v^+, \theta_w^+, \theta_w^-). \quad (9)$$

**Regularisation loss.** The final component of our loss $l$ consists of regularisation loss terms which directly encourage the alignment of the audio and visual embeddings with the text embeddings while preserving the information from their respective input modality. For this, we add two reconstruction blocks $A_{rec}$ and $V_{rec}$, such that $\phi_a^{rec} = A_{rec}(\theta_a)$ and $\phi_v^{rec} = V_{rec}(\theta_v)$, $\phi_a^{rec}, \phi_v^{rec} \in \mathbb{R}^{k_f}$. $A_{rec}$ and $V_{rec}$ each consist of a linear projection layer followed by batch normalisation, ReLU, and dropout with dropout rate $r_{dec}$:

$$l_r = d(\phi_v^{rec}, \phi_v) + d(\phi_a^{rec}, \phi_a) \\ + d(\theta_v, \theta_w) + d(\theta_a, \theta_w). \quad (10)$$

## 4. Experiments

We apply our AVCA model to audio-visual GZSL and ZSL for video classification. In this section, we first describe our proposed benchmark (Section 4.1). We discuss implementation details (Section 4.2), and then ablate the choice of different model components and loss functions (Section 4.5). Finally, we compare AVCA to state-of-the-art baseline methods for (G)ZSL (Section 4.3), and provide a detailed qualitative analysis of the learnt multi-modal embeddings (Section 4.4).

| Dataset | # classes | | # videos | | | | |
|---|---|---|---|---|---|---|---|
| | all | tr / v(U) / ts(U) | tr | v (S) | v (U) | ts (S) | ts (U) |
| VGGSound-GZSL | 276 | 138 / 69 / 69 | 70351 | 7817 | 3102 | 9032 | 3450 |
| UCF-GZSL | 51 | 30 / 12 / 9 | 3174 | 353 | 1467 | 555 | 1267 |
| ActivityNet-GZSL | 200 | 99 / 51 / 50 | 9204 | 1023 | 4307 | 1615 | 4199 |

Table 1. Statistics for our VGGSound, UCF, and ActivityNet (G)ZSL datasets, showing the number (#) of classes and videos in our splits (tr: train, v: validation, ts: test; S: seen, U: unseen).

## 4.1. Audio-Visual GZSL Benchmark

In this section, we propose three benchmark datasets for audio-visual GZSL curated from the VGGSound [19], UCF101 [64], and ActivityNet [26] datasets (summarised in Table 1)[1], and introduce our training and evaluation protocol.

**Dataset statistics.** For our proposed audio-visual GZSL splits, we include classes contained in the Sports1M [37] dataset only in our seen subsets to allow the use of feature extractors pretrained on Sports1M without leakage of information to unseen classes.

Our GZSL splits for the three datasets consist of a training set (tr), a validation set which is divided into a subset with samples from seen classes (v(S)) and another one with unseen classes (v(U)). Finally, we provide a test set consisting of seen classes (ts(S)) and unseen classes (ts(U)). The training set and the seen validation subset share the same classes with a ratio of 0.9/0.1 with respect to the number of videos. The subsets {tr ∪ v(U) ∪ v(S)} and ts(S) share the same classes and were split to have a ratio of 0.9/0.1 with respect to the number of videos.

*VGGSound* [19] is a large audio-visual dataset with 309 classes and over 200k videos. The videos can be grouped into the 9 categories *animals, home, music, nature, people, sports, tools, vehicle*, and *others*. For our VGGSound-GZSL split, we exclude videos from the *others* category and all samples from v(U) and ts(U) that were used to train SeLaVi [10], resulting in 93,752 videos in 276 classes. The 42 classes that overlap with the Sports1M dataset are only used as training classes for GZSL.

*UCF101* [64] is a video action recognition dataset which consists of over 13k videos in 101 classes. We use the subset of UCF101 which contains audio information. This results in a total of 6,816 videos for 51 classes. Previous (visual-only) methods repeatedly split the dataset into random seen and unseen classes. The 30 classes contained in the Sports1M dataset are not selected as unseen classes.

*ActivityNet* [26] is an action recognition dataset with 20k videos in 200 classes of varying duration. Again, we propose the ActivityNet-GZSL split ensuring that the 99

---

[1]VGGSound is covered by a Creative Commons license: https://creativecommons.org/licenses/by/4.0/, ActivityNet by the MIT license: https://github.com/activitynet/ActivityNet/blob/master/LICENSE.

classes contained in the Sports1M dataset are not selected as unseen classes.

**Training and evaluation protocol.** We introduce a unified training and evaluation protocol for our GZSL benchmarks. We follow this protocol to train and test all models, including AVCA and the baselines that we compare to.

We propose a two-stage training and evaluation protocol for GZSL. In the first stage, we train the models on the training set (tr), using the subsets of seen validation classes (v(S)) and unseen validation classes (v(U)) to determine the GZSL parameters, for instance for calibrated stacking [16].

In the second training stage, we re-train the models on the training (tr) and full validation set {v(S) ∪ v(U)} using the GZSL parameters determined during the first training stage. Our final models are then evaluated on the test set {ts(S) ∪ ts(U)}. ts(S) contains samples from the same classes as the training classes with no overlap between training samples for the second stage and the test samples. In particular, there is no class overlap between v(U) and ts(U).

**Evaluation metrics.** Following [78], we propose to evaluate all models using the mean class accuracy. For GZSL, we evaluate the models on the full test set {ts(S) ∪ ts(U)}, and report the averaged performance on the unseen (U) and seen (S) classes. Furthermore, we compute their harmonic mean $HM = \frac{2US}{U+S}$. We report the ZSL performance by evaluating only on the subset ts(U).

## 4.2. Experimental Setting

For each video, we use the self-supervised SeLaVi [10] framework pretrained on VGGSound [19] to extract audio and visual features for each second in a video. In our VGGSound-GZSL split, there is no overlap between videos in the unseen test and unseen validation sets and videos that were used for pre-training SeLaVi. We average the per-second features extracted using SeLaVi prior to the two-layer MLP heads to obtain 512-dimensional per-video audio and visual features. We provide additional results for using features extracted from audio and video classification networks in the supplementary material.

All networks were optimised for GZSL performance (HM) and we do not train separate networks for GZSL and ZSL. The training for the first stage was done for 50 epochs. We selected the number of training epochs for the second stage based on the GZSL performance on the validation set in the first stage. To eliminate the bias that the ZSL methods have towards seen classes, we used calibrated stacking [16] on the interval $[0, 3]$ with a step size of 0.2. For AVCA, $k_{input}$ was set to 512 and the size of the word2vec embedding, $k_{w2v}$, was set to 300. We used dropout rates $r_{dec}/r_{enc}/r_{proj}$ of 0.5/0.2/0.3 for UCF-GZSL, 0.1/0.2/0.2 for Activity-GZSL, and 0.1/0/0 for VGGSound-GZSL. The layer dimensions were set to $k_f = 300$, $k_{fhidd} = 512$, $k_{attnhidd} = 64$, and $k_{proj} = 64$. We used 3 heads for

| Method type | Model | VGGSound-GZSL | | | | UCF-GZSL | | | | ActivityNet-GZSL | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | S | U | HM | ZSL | S | U | HM | ZSL | S | U | HM | ZSL |
| ZSL | ALE [6] | 0.28 | 5.48 | 0.53 | 5.48 | 57.59 | 14.89 | 23.66 | 16.32 | 2.63 | 7.87 | 3.94 | 7.90 |
| | SJE [7] | 48.33 | 1.10 | 2.15 | 4.06 | 63.10 | 16.77 | 26.50 | 18.93 | 4.61 | 7.04 | 5.57 | 7.08 |
| | DEVISE [27] | 36.22 | 1.07 | 2.08 | 5.59 | 55.59 | 14.94 | 23.56 | 16.09 | 3.45 | 8.53 | 4.91 | 8.53 |
| | APN [85] | 7.48 | 3.88 | 5.11 | 4.49 | 28.46 | 16.16 | 20.61 | 16.44 | 9.84 | 5.76 | 7.27 | 6.34 |
| | f-VAEGAN-D2 [80] | 12.77 | 0.95 | 1.77 | 1.91 | 17.29 | 8.47 | 11.37 | 11.11 | 4.36 | 2.14 | 2.87 | 2.40 |
| Audio-visual ZSL | CJME [56] | 8.69 | 4.78 | 6.17 | 5.16 | 26.04 | 8.21 | 12.48 | 8.29 | 5.55 | 4.75 | 5.12 | 5.84 |
| | AVGZSLNet [44] | 18.05 | 3.48 | 5.83 | 5.28 | 52.52 | 10.90 | 18.05 | 13.65 | 8.93 | 5.04 | 6.44 | 5.40 |
| | AVCA | 14.90 | 4.00 | **6.31** | **6.00** | 51.53 | 18.43 | **27.15** | **20.01** | 24.86 | 8.02 | **12.13** | **9.13** |

Table 2. Evaluating our AVCA model and state-of-the-art audio-visual ZSL methods and adapted ZSL methods for GZSL and ZSL on the VGGSound, UCF, and ActivityNet (G)ZSL benchmarks. We report the mean class accuracy on the seen (S) and unseen (U) test classes, and their harmonic mean (HM) for GZSL performance. The ZSL performance is evaluated on the test subset from unseen classes.

self-attention. The loss margin hyperparameter, $\mu$, was set to 1. We used a batchsize of 256 for UCF-GZSL and ActivityNet-GZSL, and 64 for VGGSound-GZSL. We used the Adam optimiser [39] with an initial learning rate of 0.001 which was reduced by a factor of 0.1 when the GZSL performance plateaued with a patience of 3 epochs.

### 4.3. Comparing with the State of the Art

**Compared methods.** In our benchmark study, we include four image-based state-of-the-art methods and one generative method for (G)ZSL which we adapt to take audio-visual features as inputs. For this, we concatenate the audio and visual features and use those as inputs instead of image features. Moreover, we compare to current state-of-the-art methods for audio-visual GZSL [44, 56]. Here, we describe each of the methods that we compare to in more detail.

**ALE** [6] learns a linear mapping between the input features and the ground-truth embeddings, such that the projection of the input features is close to the ground-truth embedding for the corresponding class. For this, it uses a weighted approximate ranking objective [72]. **SJE** [7] computes the dot product between linearly mapped input features and the ground-truth embedding of all negative classes. The highest dot product for each example is chosen and then minimised. **DEVISE** [27] also computes the dot product between the output of a linear projection and the negative class embeddings and it minimises the sum of these dot products. **APN** [85] is the current non-generative state-of-the-art method for image-based ZSL. APN is based on the assumption that the ground-truth embeddings contain visual class attributes. Prototypes are used to map the attributes from the ground-truth embeddings to relevant locations in the image. **f-VAEGAN-D2** [80] is a generative ZSL method which learns to generate synthetic features for unseen classes. Then, a classifier is trained on real examples from seen classes and synthetic examples from

unseen classes. **CJME** [56] proposed the task of audio-visual GZSL for video classification on the AudioSetZSL dataset. It embeds audio, video and text into a joint embedding space and uses proximity in the embedding space to select the classification output at test time. **AVGZSLNet** [44] builds on [56] and is the current state-of-the-art method for audio-visual GZSL for video classification. One of the main strengths of this method is its use of triplet losses to leverage information from negative examples.

**Results.** We compare our AVCA framework to recent methods for (G)ZSL in Table 2 on the VGGSound-GZSL, UCF-GZSL, and ActivityNet-GZSL datasets. AVCA obtains the best results on all three datasets. On VGGSound-GZSL, AVCA obtains a HM of 6.31% for GZSL and a ZSL performance of 6.00% compared to 6.17% HM for CJME and a ZSL performance of 5.59% for DEVISE. On the UCF-GZSL dataset, our AVCA model outperforms SJE for GZSL with a performance of 27.15% compared to 26.50%, and we obtain a stronger ZSL performance of 20.01% compared to 18.93%. On ActivityNet-GZSL, AVCA outperforms APN, with a GZSL performance of 12.13% compared to 7.27%. For ZSL, AVCA is stronger than DEVISE with a score of 9.13% compared to 8.53%. It can be observed that in some cases U is higher than S. This is due to the use of calibrated stacking [16] as described in [46].

### 4.4. Qualitative Results

We present a qualitative analysis of the learnt multi-modal embeddings in Fig. 3. The t-SNE visualisations [73] for a subset of ActivityNet-GZSL classes show the differences between the audio and visual input features and the learnt multi-modal embeddings. We provide additional qualitative results for VGGSound-GZSL and UCF-GZSL in the supplementary material. It can be seen in Fig. 3a that the input audio features are not as well-separated and clustered as the visual features shown in Fig. 3b. However, the vi-
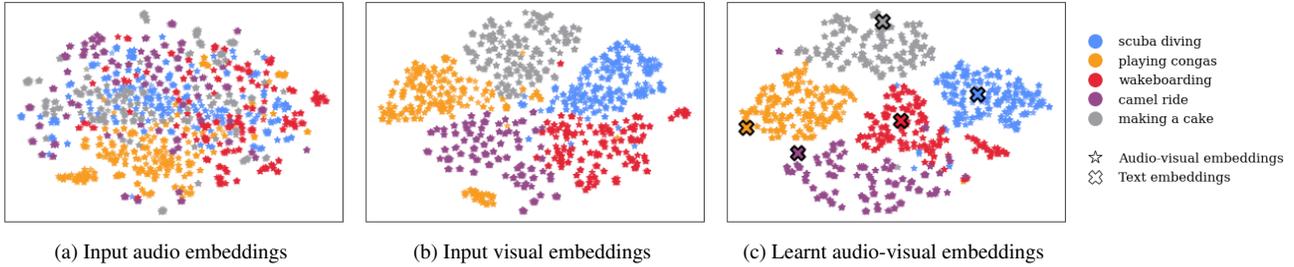
| (a) Input audio embeddings | (b) Input visual embeddings | (c) Learnt audio-visual embeddings |

Figure 3. t-SNE visualisation for three seen (*scuba diving, playing congas, wakeboarding*) and two unseen (*camel ride, making a cake*) test classes from ActivityNet-GZSL, showing embeddings extracted with SeLaVi [10] for (a) audio and (b) visual features. (c) Learnt audio-visual embeddings of our model. Projected textual class label embeddings are visualised with a cross with black boundary.

| Model | VGGSound-GZSL | | UCF-GZSL | | ActivityNet-GZSL | |
|---|---|---|---|---|---|---|
| | HM | ZSL | HM | ZSL | HM | ZSL |
| Visual branch | 4.83 | 4.06 | 20.92 | 14.16 | 7.53 | 6.49 |
| Audio branch | 3.84 | 3.83 | 11.78 | 10.78 | 4.19 | 4.06 |
| AVCA | **6.31** | **6.00** | **27.15** | **20.01** | **12.13** | **9.13** |

Table 3. Influence of *training* AVCA with different modalities for GZSL and ZSL on the VGGSound-GZSL, UCF-GZSL and ActivityNet-GZSL datasets measuring the harmonic mean (HM) for GZSL and the mean class accuracy for ZSL. Using both modalities yields the strongest GZSL and ZSL performances.

| Model | VGGSound-GZSL | | UCF-GZSL | | ActivityNet-GZSL | |
|---|---|---|---|---|---|---|
| | HM | ZSL | HM | ZSL | HM | ZSL |
| W/o x-att | 6.02 | 4.81 | 26.82 | 18.37 | 6.50 | 5.64 |
| Visual with x-att | **6.63** | 4.78 | 27.11 | 17.22 | 9.50 | 6.89 |
| Audio with x-att | 4.93 | 5.01 | 18.61 | 16.05 | 11.05 | 8.78 |
| AVCA | 6.31 | **6.00** | **27.15** | **20.01** | **12.13** | **9.13** |

Table 4. Using different components of AVCA for GZSL and ZSL on VGGSound-GZSL, UCF-GZSL and ActivityNet-GZSL. Audio (Visual) with x-att uses the visual (audio) modality only for the cross-attention. W/o x-att optimises each branch in isolation and their output predictions are averaged. x-att denotes cross-attention.

sual features also contain classes, such as *playing congas* and *scuba diving*, which are not clustered cleanly. It can be observed in Fig. 3c that our model produces multi-modal features that improve over the clustering of the input embeddings for both, seen and unseen classes. For instance, the cluster separation between the seen class *playing congas* and the unseen class *making a cake* improves significantly, even though the unseen class is not used for training.

## 4.5. Ablation Analysis

Here, we analyse how different architectural choices and loss components for AVCA impact the performances on VGGSound-GZSL, ActivityNet-GZSL, and UCF-GZSL.

**Evaluating different modalities.** In Table 3, we compare our multi-modal AVCA model to training our architecture with only unimodal inputs. In this case, we remove the cross-modal attention block and train each unimodal branch in isolation. The visual branch obtains a better performance than the audio branch with a GZSL performance (HM) of 7.53% vs. 4.19% on the ActivityNet-GZSL dataset. A similar pattern can be observed for the ZSL performance with 6.49% vs. 4.06% for the visual and audio branch respectively. This trend is also exhibited on the UCF-GZSL and VGGSound-GZSL datasets, suggesting that the visual input features provide richer information about the video content than the audio inputs. Nevertheless, jointly training AVCA with both input modalities gives significant improvements over using each of them individually with a GZSL performance of 12.13% and a ZSL performance of 9.13% on the

ActivityNet-GZSL dataset. This confirms that the complementary information from the audio and visual inputs is highly beneficial for GZSL and ZSL for video classification. We provide the S/U performances for Table 3 in the supplementary material.

**Evaluating the cross-modal attention block.** Next, we investigate the effect of using our cross-modal attention block in Table 4. To obtain results without using cross-attention (W/o x-att), each branch is optimised individually. For evaluation, we compute the distances between the outputs of both branches and $\theta_w$ for each class, and then average the distances computed by both branches. The GZSL and ZSL performances drop dramatically when not using the cross-attention block from 12.13% and 9.13% for AVCA to 6.50% and 5.64% for GZSL and ZSL scores respectively on the ActivityNet-GZSL dataset. The pattern is similar for VGGSound-GZSL and UCF-GZSL, confirming the importance of our cross-modal attention block for sharing information between the input modalities.

Furthermore, we compare optimising our full AVCA model to using only the visual (Visual with x-att) or only the audio branch (Audio with x-att) for training. Using only the visual branch entails removing $A_{rec}$ and $A_{proj}$ along with their associated losses from the audio branch but keeping the cross-attention. This experiment is repeated for the audio branch by removing the corresponding components from the visual branch. Jointly optimising both branches provides better results than using only

| Model output | VGGSound-GZSL | | UCF-GZSL | | ActivityNet-GZSL | |
|---|---|---|---|---|---|---|
| | HM | ZSL | HM | ZSL | HM | ZSL |
| AVCA ($\theta_a$) | 5.18 | 4.87 | 25.98 | 18.25 | **12.54** | **9.23** |
| AVCA ($\theta_v$) | **6.31** | **6.00** | **27.15** | **20.01** | 12.13 | 9.13 |
| AVCA ($\theta_a, \theta_v$) | 5.90 | 5.42 | 25.78 | 19.30 | 12.17 | 8.95 |
| AVCA ($min(\theta_a, \theta_v)$) | 6.10 | 5.36 | 25.86 | 18.39 | 12.45 | 9.08 |

Table 5. Influence of using the outputs of the audio and visual branches $\theta_a$ and $\theta_v$ separately, or using both jointly ($\theta_a, \theta_v$) for *evaluation* on VGGSound-GZSL, UCF-GZSL and ActivityNet-GZSL. All models were trained with $\theta_a$ and $\theta_v$.

| Model | VGGSound-GZSL | | UCF-GZSL | | ActivityNet-GZSL | |
|---|---|---|---|---|---|---|
| | HM | ZSL | HM | ZSL | HM | ZSL |
| $l-l_t$ | 5.06 | 4.84 | 18.51 | 19.17 | 8.39 | **9.54** |
| $l-l_{rec}$ | 5.92 | 5.22 | 24.32 | 17.20 | 9.59 | 6.93 |
| $l-l_{ct}$ | 6.31 | 4.87 | 17.88 | 17.51 | 11.20 | 8.99 |
| $l-l_w$ | 5.18 | 4.93 | 20.75 | 16.41 | 9.08 | 8.00 |
| $l-l_r$ | 6.24 | 4.43 | 21.31 | 14.02 | 11.14 | 7.94 |
| $l$ | **6.31** | **6.00** | **27.15** | **20.01** | **12.13** | 9.13 |

Table 6. Comparing training AVCA with our full loss function $l$ to removing individual components $l_t$, $l_{rec}$, $l_{ct}$, $l_w$, or $l_r$, on the GZSL and ZSL performance on the VGGSound-GZSL, UCF-GZSL and ActivityNet-GZSL datasets.

one of the branches on ActivityNet-GZSL and UCF-GZSL. On ActivityNet-GZSL, we obtain a GZSL performances of 12.13% compared to 11.05% and 9.50% for using only the audio and visual branches respectively. Interestingly, for the VGGSound-GZSL dataset, the Visual with x-att model yields a slightly stronger GZSL performance than our full AVCA model, with a HM of 6.63% compared to 6.31%. This is in line with the Audio branch performing worse than the Visual branch on VGGSound-GZSL (Table 3). However, the joint optimisation of AVCA gives the best results.

**Evaluating different modalities as output.** In Table 5, we investigate the effect of evaluating our full trained AVCA model using only the output features from the audio ($\theta_a$) or the visual ($\theta_v$) branch, or from both branches together (($\theta_a, \theta_v$) and $min(\theta_a, \theta_v)$). For AVCA($\theta_a, \theta_v$), we compute the distance $|\theta_a - \theta_w|_2 + |\theta_v - \theta_w|_2$. AVCA($min(\theta_a, \theta_v)$) uses the embedding from the modality that has the smallest distance to a word embedding. The class corresponding to the closest text embedding resembles the class prediction.

Using the visual branch gives the strongest performance on VGGSound-GZSL/UCF-GZSL with a HM of 6.31%/27.15% vs 5.18%/25.98% for the audio branch. On ActivityNet-GZSL, the audio branch yields slightly better results (HM of 12.54% vs. 12.13% for the visual branch). Both AVCA($\theta_a, \theta_v$) and AVCA($min(\theta_a, \theta_v)$) obtain lower scores that $\theta_v$. The best results (highest averaged HM) across all three datasets are produced when using the visual branch only. However, as the cross-attention block fuses the audio and visual modalities, both branches contain multimodal information from both input modalities.

**Evaluating different loss functions.** Finally, we analyse the impact of using different loss functions for training AVCA on the GZSL and ZSL performance in Table 6. We observe that using our full loss $l$ provides the strongest GZSL results (HM) on the UCF-GZSL, VGGSound-GZSL, and ActivityNet-GZSL datasets by a large margin. On ActivityNet-GZSL, omitting $l_t$ for training our model $(l-l_t)$ provides slightly stronger ZSL results than using our full loss $l$ with a mean class accuracy of 9.54% compared to 9.13%. However, the GZSL performance is significantly better when using $l$ with a HM of 12.13% compared to 8.39% when using $l - l_t$. Our loss ablations confirm that

our strong overall performance on all three datasets is only obtained when training with our full proposed loss function.

### 4.6. Limitations and Discussion

Our proposed GZSL benchmark datasets pose an extremely challenging setting, since the underlying datasets span a wide variety of classes (*e.g.* including *wakeboarding* and *making a cake* for the ActivityNet dataset). Our AVCA leverages the varied audio-visual input information effectively, resulting in more robust GZSL performance than the related methods. However, AVCA uses temporally averaged audio-visual input information, and hence does not consider fine semantic details. Furthermore, our model relies on multi-modal input data and cannot be used when only one modality is available.

## 5. Conclusion

We introduced three new benchmarks for audio-visual (generalised) zero-shot learning for video classification on the VGGSound, UCF, and ActivityNet datasets. We proposed a framework for (G)ZSL from audio-visual data which learns to align the audio-visual embeddings with textual label embeddings. Furthermore, we provided baseline performances for seven (G)ZSL methods, and show that our model outperforms them for GZSL and ZSL on our new benchmarks. Finally, we provided a qualitative analysis of the learnt multi-modal embeddings. We hope that our proposed benchmarks will enable and encourage further research into audio-visual zero-shot learning.

# Supplementary Material: Audio-visual Generalised Zero-shot Learning with Cross-modal Attention and Language

In this supplementary material, we include additional qualitative results (Appendix A) and quantitative results (Appendix B) for our proposed audio-visual (G)ZSL framework.

## A. Additional Qualitative Results

We provide additional qualitative results for our proposed AVCA model for the tasks of audio-visual GZSL and ZSL. We present t-SNE visualisations for the learnt audio-visual embeddings on the VGGSound-GZSL and UCF-GZSL datasets in Fig. 4 and Fig. 5.

In Fig. 4a, we can observe that the input audio features do not demonstrate a clear separation between the visualised classes for the VGGSound-GZSL dataset. The visual features exhibit a better clustering as can been seen in Fig. 4b. However, the visual features also include classes, such as *elephant trumpeting* and *wood thrush calling*, that are not clustered cleanly. Our AVCA model outputs multimodal features that improve the clustering for both, seen and unseen classes (Fig. 4c). The learnt features for the two unseen classes *elephant trumpeting* and *wood thrush calling* are clustered and well-separated as opposed to the input features. This is impressive, since both classes were not included in the training set.

Similarly, for the UCF-GZSL dataset, we can observe in Fig. 5a that the input audio features are not grouped ac-
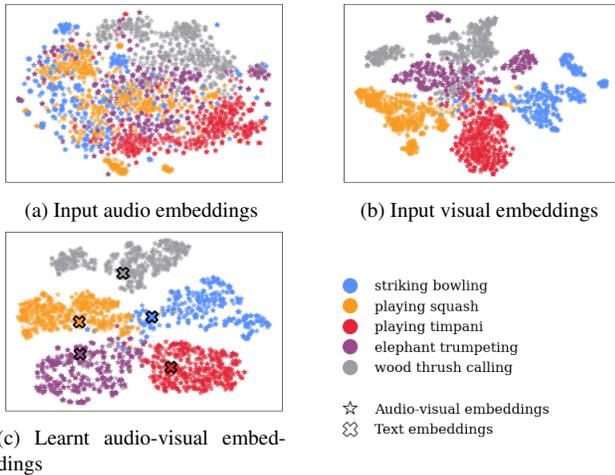


(a) Input audio embeddings

(b) Input visual embeddings



(c) Learnt audio-visual embeddings

- ● baby crawling
- ● basketball dunk
- ● bowling
- ● band marching
- ● playing flute

- ☆ Audio-visual embeddings
- ⊗ Text embeddings

Figure 5. t-SNE visualisation for three seen (*baby crawling, basketball dunk, bowling*) and two unseen (*band marching, playing flute*) test classes from the UCF-GZSL dataset, showing (a) audio and (b) visual features extracted with SeLaVi [10], and (c) learnt audio-visual embeddings of our model. Textual class label embeddings are visualised with a cross.
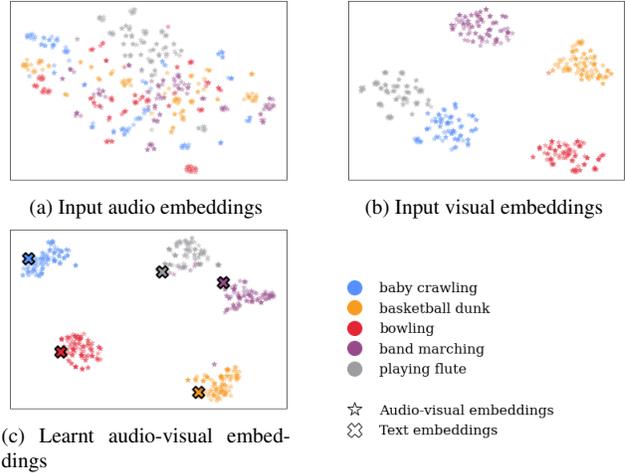
cording to classes. In contrast, the visual input embeddings mostly exhibit a clear clustering of different classes. However, the classes *baby crawling* and *playing flute* are not well-separated as can be seen in Fig. 5b. This improves through learning, since the learnt audio-visual features in Fig. 5c show a clear divide between those two classes. In addition to that, the output embeddings for the unseen classes *band marching* and *playing flute* are overwhelmingly clustered well, too.

To summarise, our model learns to cluster both seen and unseen classes for different datasets by transferring information from the training data to unseen classes at test time.

## B. Additional Quantitative Results

In this section, we provide additional quantitative results obtained with our AVCA. We present results for training and evaluating our AVCA model with a different set of input features in Appendix B.1. In particular, we use features extracted from networks that were pretrained for audio and video classification. We perform an additional ablation study that gradually transforms AVCA into AVGZSLNet [44] in Appendix B.2. Complete results that include the U and S performance for Table 3 in the main paper are provided in Appendix B.3. Finally, we give details about the number of parameters and GFLOPS required for training our AVCA model in Appendix B.4



(a) Input audio embeddings

(b) Input visual embeddings



(c) Learnt audio-visual embeddings

- ● striking bowling
- ● playing squash
- ● playing timpani
- ● elephant trumpeting
- ● wood thrush calling

- ☆ Audio-visual embeddings
- ⊗ Text embeddings

Figure 4. t-SNE visualisation for three seen (*striking bowling, playing squash, playing timpani*) and two unseen (*elephant trumpeting, wood thrush calling*) test classes from the VGGSound-GZSL dataset, showing (a) audio and (b) visual features extracted with SeLaVi [10], and (c) learnt audio-visual embeddings of our model. Textual class label embeddings are visualised with a cross.

| Method type | Model | VGGSound-GZSL$^{cls}$ | | | | UCF-GZSL$^{cls}$ | | | | ActivityNet-GZSL$^{cls}$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | S | U | HM | ZSL | S | U | HM | ZSL | S | U | HM | ZSL |
| ZSL | ALE [6] | 26.13 | 1.72 | 3.23 | 4.97 | 45.42 | 29.09 | 35.47 | 32.30 | 0.89 | 6.16 | 1.55 | 6.16 |
| | SJE [7] | 16.94 | 2.72 | 4.69 | 3.22 | 19.39 | 32.47 | 24.28 | 32.47 | 37.92 | 1.22 | 2.35 | 4.35 |
| | DEVISE [27] | 29.96 | 1.94 | 3.64 | 4.72 | 29.58 | 34.80 | 31.98 | 35.48 | 0.17 | 5.84 | 0.33 | 5.84 |
| | APN [85] | 6.46 | 6.13 | 6.29 | 6.50 | 13.54 | 28.44 | 18.35 | 29.69 | 3.79 | 3.39 | 3.58 | 3.97 |
| Audio-visual ZSL | CJME [56] | 10.86 | 2.22 | 3.68 | 3.72 | 33.89 | 24.82 | 28.65 | 29.01 | 10.75 | 5.55 | 7.32 | 6.29 |
| | AVGZSLNet [44] | 15.02 | 3.19 | 5.26 | 4.81 | 74.79 | 24.15 | 36.51 | 31.51 | 13.70 | 5.96 | 8.30 | 6.39 |
| | AVCA | 12.63 | 6.19 | **8.31** | **6.91** | 63.15 | 30.72 | **41.34** | **37.72** | 16.77 | 7.04 | **9.92** | **7.58** |

Table 7. Evaluating AVCA and state-of-the-art (G)ZSL methods for audio-visual GZSL and ZSL on the VGGSound, UCF, and ActivityNet (G)ZSL$^{cls}$ benchmarks using features extracted from audio/video classification networks. We report the mean class accuracy on the seen (S) and unseen (U) test classes, and their harmonic mean (HM) for GZSL performance. The ZSL performance is evaluated on the test subset of samples from unseen classes.

| Dataset | # classes | | | | # videos |
|---|---|---|---|---|---|
| | all | tr | v(U) | ts(U) | ts(U) |
| VGGSound-GZSL$^{cls}$ | 271 | 138 | 69 | 64 | 3200 |
| UCF-GZSL$^{cls}$ | 48 | 30 | 12 | 6 | 845 |
| ActivityNet-GZSL$^{cls}$ | 198 | 99 | 51 | 48 | 4052 |

Table 8. Statistics for our VGGSound, UCF, and ActivityNet (G)ZSL$^{cls}$ datasets, showing the number (#) of classes in our splits (tr: train, v: validation, ts: test; S: seen, U: unseen). $^{cls}$ indicates the dataset splits that allow to use VGGish features pretrained on YouTube-8M. The full details about the dataset splits can be found at https://github.com/ExplainableML/AVCA-GZSL.

## B.1. Using features extracted audio/video classification networks

We additionally trained and tested our model and the baseline models using features extracted from audio and video classification networks (instead of the SeLaVi [10] features used in the main paper). In particular, the visual features were extracted with C3D [70], pretrained for video classification on Sports1M [37]. The audio features were extracted with VGGish [34], pretrained for audio classification on Youtube-8M [1]. We averaged the extracted features across time, resulting in a 4096-dimensional visual feature and a 128-dimensional audio feature for each video.

However, to use the audio features extracted from a network that was pretrained on Youtube-8M, we removed the test unseen classes from the VGGSound-GZSL, UCF-GZSL, and ActivityNet-GZSL datasets that had an overlap with Youtube-8M. This resulted in slightly different dataset splits (VGGSound-GZSL$^{cls}$, UCF-GZSL$^{cls}$, and ActivityNet-GZSL$^{cls}$) detailed in Table 8.

We provide results for training and evaluating our AVCA

| Model | VGGSound-GZSL | | UCF-GZSL | | ActivityNet-GZSL | |
|---|---|---|---|---|---|---|
| | HM | ZSL | HM | ZLS | HM | ZSL |
| AVGZSLNet [44] | 5.83 | 5.28 | 18.05 | 13.65 | 6.44 | 5.40 |
| W/o x-att | 6.02 | 4.81 | 26.82 | 18.37 | 6.50 | 5.64 |
| W x-att with $l_c$ loss | 4.88 | 4.55 | 19.38 | 12.95 | 11.58 | 8.40 |
| AVCA | **6.31** | **6.00** | **27.15** | **20.01** | **12.13** | **9.13** |

Table 9. Ablation that gradually transforms our AVCA model into AVGZSLNet [44]. W/o x-att optimises each branch in isolation and their output predictions are averaged. x-att denotes cross-attention. $l_c$ loss is the loss function used to train AVGZSLNet.

and the baselines using audio and video classification features in Table 7. AVCA outperforms all the baselines on all three datasets. On VGGSound-GZSL$^{cls}$, ACVA obtains a HM of 8.31% and ZSL of 6.91% compared to a HM of 6.29% for APN and a ZSL performance of 6.50% for APN. On UCF-GZSL$^{cls}$, AVCA obtains a HM of 41.34% and a ZSL of 37.72% compared to a HM of 36.51% for AVGZSLNet and a ZSL performance of 35.48% for DEVISE. On ActivityNet-GZSL$^{cls}$, AVCA outperforms AVGZSLNet with a HM of 9.92% compared to 8.30% and a ZSL of 7.58% compared to 6.39% for AVGZSLNet. These results show that AVCA outperforms the other competitors also when using audio and video classification features, proving again that our cross-attention mechanism and training objective provide a boost in performance.

## B.2. Ablating AVCA in relation to AVGZSLNet

We additionally perform an ablation study that gradually transforms the AVCA model into AVGZSLNet [44] in Table 9. We show how our model components influence the (G)ZSL performance, resulting in our AVCA model that outperforms AVGZSLNet on all three datasets. For this ablation, we use the SeLaVi [10] features and the same setup as in the main paper. W/o x-att corre-

| Model | VGGSound-GZSL | | | UCF-GZSL | | | ActivityNet-GZSL | | |
|---|---|---|---|---|---|---|---|---|---|
| | S | U | HM | S | U | HM | S | U | HM |
| Visual branch | 7.02 | 3.68 | 4.83 | 50.18 | 13.21 | 20.92 | 11.80 | 5.53 | 7.53 |
| Audio branch | 7.74 | 2.55 | 3.84 | 12.99 | 10.78 | 11.78 | 4.56 | 3.87 | 4.19 |
| AVCA | 14.90 | 4.00 | **6.31** | 51.53 | 18.43 | **27.15** | 24.86 | 8.02 | **12.13** |

Table 10. Influence of *training* AVCA with different modalities for GZSL on the VGGSound-GZSL, UCF-GZSL and ActivityNet-GZSL datasets measuring the GZSL performance on seen (S) and unseen (U) test classes and their harmonic mean (HM). Using both modalities yields the strongest GZSL performances.

sponds to AVGZSLNet trained with our loss function (without our cross-attention). It can be observed that W/o x-att provides improvements on UCF-GZSL, with a HM of 26.82% compared to 18.05% and a ZSL performance of 18.37% compared to 13.65%. W x-att with $l_c$ loss corresponds to AVGZSLNet with cross-attention and with the loss function proposed for AVGZSLNet. In this case, it can be observed that the cross-attention improves the results over AVGZSLNet with a HM of 11.58% compared to 6.44% and ZSL performance of 8.40% compared to 5.40% on ActivityNet-GZSL. These improvements can also be observed on the other datasets, showing that our novel loss and our cross-attention mechanism improve the performance over AVGZSLNet.

## B.3. Extended results for training AVCA with different modalities

In this section, we extend the ablation study that uses different modalities for training (Table 3 in the main paper) by adding the performance on the seen (S) and unseen (U) test classes for all the datasets in Table 10.

On all three datasets it can be observed that there is an increase in both seen and unseen performance when using AVCA compared to using the Visual branch or the Audio branch. On VGGSound-GZSL, we can observe that the S performance for AVCA is 14.90% compared to 7.74% for the Visual branch. The U performance on VGGSound-GZSL is also stronger for AVCA than for the Visual branch, with a score of 4.00% compared to 3.68%. On the UCF-GZSL dataset, the S performance increases only slightly, from 50.18% for the Visual branch to 51.53% for AVCA. However, there is a significant increase in the U performance, from 13.21% for the Visual branch to 18.43% for AVCA. Finally, on ActivityNet-GZSL, AVCA yields a S score of 24.86% compared to 11.80% for the Visual branch. The U performance increases from 5.53% for the Visual branch to 8.02% for AVCA. These results show that the S/U performance increases significantly when using AVCA compared to the Visual branch or the Audio branch, leading to better HM/ZSL performances.

## B.4. Number of parameters in AVCA.

AVCA contains 1.69M parameters in total, which is comparable to the 1.32M parameters used in AVGZSLNet [44]. ALE/SJE/DEVISE are significantly smaller with only 307.2k parameters. AVCA has a computational complexity of 2.36 GFLOPS, while AVGZSLNet has a computational complexity of 1.38 GFLOPS. Again, the fewest GFLOPS are required for ALE/SJE/DEVISE which have a computational complexity of 0.32 GFLOPS. These statistics show that AVCA is comparable to AVGZSLNet while providing significantly better results on all three datasets.

## References

[1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016. 10

[2] Triantafyllos Afouras, Yuki M Asano, Francois Fagan, Andrea Vedaldi, and Florian Metze. Self-supervised object detection from audio-visual correspondence. In *ECCV*, 2020. 2

[3] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. *IEEE TPAMI*, 2018. 2

[4] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Asr is all you need: Cross-modal distillation for lip reading. In *ICASSP*, 2020. 2

[5] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In *ECCV*, 2020. 2

[6] Zeynep Akata, Florent Perronnin, Zaid Harchaoui, and Cordelia Schmid. Label-embedding for image classification. *IEEE TPAMI*, 2015. 1, 2, 6, 10

[7] Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output embeddings for fine-grained image classification. In *CVPR*, 2015. 1, 2, 6, 10

[8] Humam Alwassel, Dhruv Mahajan, Lorenzo Torresani, Bernard Ghanem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. In *NeurIPS*, 2020. 2

[9] Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *ECCV*, 2018. 2

[10] Yuki M. Asano, Mandela Patrick, Christian Rupprecht, and Andrea Vedaldi. Labelling unlabelled videos from scratch with multi-modal self-supervision. In *NeurIPS*, 2020. 2, 5, 7, 9, 10

[11] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *NeurIPS*, 2016. 2

[12] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 3

[13] Mina Bishay, Georgios Zoumpourlis, and Ioannis Patras. Tarn: Temporal attentive relation network for few-shot and zero-shot action recognition. In *BMVC*, 2019. 1, 2

[14] Biagio Brattoli, Joseph Tighe, Fedor Zhdanov, Pietro Perona, and Krzysztof Chalupka. Rethinking zero-shot video classification: End-to-end training for realistic applications. In *CVPR*, 2020. 1, 2

[15] Andrew Brown, Vicky Kalogeiton, and Andrew Zisserman. Face, body, voice: Video person-clustering with multiple modalities. In *ICCV*, 2021. 2

[16] Wei-Lun Chao, Soravit Changpinyo, Boqing Gong, and Fei Sha. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *ECCV*, 2016. 5, 6

[17] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Audio-visual synchronisation in the wild. In *BMVC*, 2021. 2

[18] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *CVPR*, 2021. 2

[19] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, 2020. 2, 5

[20] Yanbei Chen, Yongqin Xian, A. Sophia Koepke, Ying Shan, and Zeynep Akata. Distilling audio-visual knowledge by compositional contrastive learning. In *CVPR*, 2021. 2

[21] Ying Cheng, Ruize Wang, Zhihao Pan, Rui Feng, and Yuejie Zhang. Look, listen, and attend: Co-attention network for self-supervised audio-visual representation learning. In *ACM MM*, 2020. 2

[22] Jeong Choi, Jongpil Lee, Jiyoung Park, and Juhan Nam. Zero-shot learning for audio-based music classification and tagging. In *ISMIR*, 2019. 2

[23] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *ACCV*, 2016. 2

[24] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 2

[25] Joshua P Ebeneze, Yongjun Wu, Hai Wei, Sriram Sethuraman, and Zongyi Liu. Detection of audio-video synchronization errors via event detection. In *ICASSP*, 2021. 2

[26] Bernard Ghanem Fabian Caba Heilbron, Victor Escorcia and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015. 2, 5

[27] Andrea Frome, Greg Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *NeurIPS*, 2013. 2, 6, 10

[28] Chuang Gan, Deng Huang, Peihao Chen, Joshua B Tenenbaum, and Antonio Torralba. Foley music: Learning to generate music from videos. In *ECCV*, 2020. 2

[29] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019. 2

[30] Shir Goldstein and Yael Moses. Guitar music transcription from silent video. In *BMVC*, 2018. 2

[31] Shreyank N Gowda, Laura Sevilla-Lara, Kiyoon Kim, Frank Keller, and Marcus Rohrbach. A new split for evaluating true zero-shot action recognition. *arXiv preprint arXiv:2107.13029*, 2021. 2

[32] Meera Hahn, Andrew Silva, and James M Rehg. Action2vec: A crossmodal embedding approach to action learning. *arXiv preprint arXiv:1901.00484*, 2019. 1, 2

[33] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 3

[34] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *ICASSP*, 2017. 10

[35] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 3

[36] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *IJCV*, 2019. 2

[37] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. 5, 10

[38] Naji Khosravan, Shervin Ardeshir, and Rohit Puri. On attention modules for audio-visual synchronization. In *CVPR Workshop*, 2019. 2

[39] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *CoRR*, 2015. 6

[40] Elyor Kodirov, Tao Xiang, and Shaogang Gong. Semantic autoencoder for zero-shot learning. In *CVPR*, 2017. 2

[41] A Sophia Koepke, Olivia Wiles, Yael Moses, and Andrew Zisserman. Sight to sound: An end-to-end approach for visual piano transcription. In *ICASSP*, 2020. 2

[42] A. Sophia Koepke, Olivia Wiles, and Andrew Zisserman. Visual pitch estimation. In *SMC*, 2019. 2

[43] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, 2018. 2

[44] Pratik Mazumder, Pravendra Singh, Kranti Kumar Parida, and Vinay P Namboodiri. Avgzslnet: Audio-visual generalized zero-shot learning by reconstructing label features from multi-modal embeddings. In *WACV*, 2021. 1, 2, 3, 4, 6, 9, 10, 11

[45] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *ICLR*, 2013. 2

[46] Shaobo Min, Hantao Yao, Hongtao Xie, Chaoqun Wang, Zheng-Jun Zha, and Yongdong Zhang. Domain-aware visual bias eliminating for generalized zero-shot learning. In *CVPR*, 2020. 6

[47] Liliane Momeni, Triantafyllos Afouras, Themos Stafylakis, Samuel Albanie, and Andrew Zisserman. Seeing wake words: Audio-visual keyword spotting. In *BMVC*, 2020. 2

[48] Arsha Nagrani, Joon Son Chung, Samuel Albanie, and Andrew Zisserman. Disentangled speech embeddings using cross-modal self-supervision. In *ICASSP*, 2020. 2

[49] Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. In *NeurIPS*, 2021. 2

[50] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010. 3

[51] Medhini Narasimhan, Shiry Ginosar, Andrew Owens, Alexei A Efros, and Trevor Darrell. Strumming to the beat: Audio-conditioned contrastive video textures. *arXiv preprint arXiv:2104.02687*, 2021. 2

[52] Sanath Narayan, Akshita Gupta, Fahad Shahbaz Khan, Cees GM Snoek, and Ling Shao. Latent embedding feedback and discriminative features for zero-shot classification. In *ECCV*, 2020. 2

[53] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, 2018. 2

[54] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *ECCV*, 2016. 2

[55] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Learning sight from sound: Ambient sound provides supervision for visual learning. *IJCV*, 2018. 2

[56] Kranti Parida, Neeraj Matiyali, Tanaya Guha, and Gaurav Sharma. Coordinated joint multimodal embeddings for generalized audio-visual zero-shot classification and retrieval of videos. In *WACV*, 2020. 1, 2, 3, 6, 10

[57] Mandela Patrick, Yuki M Asano, Ruth Fong, João F Henriques, Geoffrey Zweig, and Andrea Vedaldi. Multi-modal self-supervision from generalized data transformations. In *NeurIPS*, 2020. 2

[58] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014. 2

[59] KR Prajwal, Liliane Momeni, Triantafyllos Afouras, and Andrew Zisserman. Visual keyword spotting with attention. In *BMVC*, 2021. 2

[60] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, 2020. 2

[61] Alina Roitberg, Manuel Martinez, Monica Haurilet, and Rainer Stiefelhagen. Towards a fair evaluation of zero-shot action recognition using external data. In *ECCV Workshop*, 2018. 2

[62] Bernardino Romera-Paredes and Philip Torr. An embarrassingly simple approach to zero-shot learning. In *ICML*, 2015. 1, 2

[63] Edgar Schonfeld, Sayna Ebrahimi, Samarth Sinha, Trevor Darrell, and Zeynep Akata. Generalized zero-and few-shot learning via aligned variational autoencoders. In *CVPR*, 2019. 1, 2

[64] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 2, 5

[65] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014. 3

[66] Kun Su, Xiulong Liu, and Eli Shlizerman. Multi-instrumentalist net: Unsupervised generation of music from body movements. *arXiv preprint arXiv:2012.03478*, 2020. 2

[67] Kun Su, Xiulong Liu, and Eli Shlizerman. How does it sound? In *NeurIPS*, 2021. 2

[68] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: weakly-supervised audio-visual video parsing. In *ECCV*, 2020. 2

[69] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *ECCV*, 2018. 2

[70] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 10

[71] Efthymios Tzinis, Scott Wisdom, Aren Jansen, Shawn Hershey, Tal Remez, Daniel PW Ellis, and John R Hershey. Into the wild with audioscope: Unsupervised audio-visual separation of on-screen sounds. In *ICLR*, 2021. 2

[72] Nicolas Usunier, David Buffoni, and Patrick Gallinari. Ranking with ordered weighted pairwise classification. In *ICML*, 2009. 6

[73] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 2008. 6

[74] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 3

[75] Vinay Kumar Verma, Gundeep Arora, Ashish Mishra, and Piyush Rai. Generalized zero-shot learning via synthesized examples. In *CVPR*, 2018. 1, 2

[76] Qian Wang and Ke Chen. Zero-shot visual recognition via bidirectional latent embedding. *IJCV*, 2017. 2

[77] Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *ECCV*, 2018. 2

[78] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE TPAMI*, 2018. 2, 5

[79] Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. Feature generating networks for zero-shot learning. In *CVPR*, 2018. 2

[80] Yongqin Xian, Saurabh Sharma, Bernt Schiele, and Zeynep Akata. f-vaegan-d2: A feature generating framework for any-shot learning. In *CVPR*, 2019. 2, 6

[81] Fanyi Xiao, Yong Jae Lee, Kristen Grauman, Jitendra Malik, and Christoph Feichtenhofer. Audiovisual slowfast networks for video recognition. *arXiv preprint arXiv:2001.08740*, 2020. 2

[82] Huang Xie, Okko Räsänen, and Tuomas Virtanen. Zero-shot audio classification with factored linear and nonlinear acoustic-semantic projections. In *ICASSP*, 2021. 2

[83] Huang Xie and Tuomas Virtanen. Zero-shot audio classification via semantic embeddings. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2021. 2

[84] Haoming Xu, Runhao Zeng, Qingyao Wu, Mingkui Tan, and Chuang Gan. Cross-modal relation-aware networks for audio-visual event localization. In *ACM MM*, 2020. 2

[85] Wenjia Xu, Yongqin Xian, Jiuniu Wang, Bernt Schiele, and Zeynep Akata. Attribute prototype network for zero-shot learning. In *NeurIPS*, 2020. 1, 2, 6, 10

[86] Wenjia Xu, Yongqin Xian, Jiuniu Wang, Bernt Schiele, and Zeynep Akata. Vgse: Visually-grounded semantic embeddings for zero-shot learning. In *CVPR*, 2022. 2

[87] Xun Xu, Timothy M Hospedales, and Shaogang Gong. Multi-task zero-shot action recognition with prioritised data augmentation. In *ECCV*, 2016. 1, 2

[88] Hanyu Xuan, Zhenyu Zhang, Shuo Chen, Jian Yang, and Yan Yan. Cross-modal attention network for temporal inconsistent audio-visual event localization. In *AAAI*, 2020. 2

[89] Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In *ICCV*, 2019. 2

[90] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018. 2

[91] Hang Zhou, Ziwei Liu, Xudong Xu, Ping Luo, and Xiaogang Wang. Vision-infused deep audio inpainting. In *ICCV*, 2019. 2

[92] Lingyu Zhu and Esa Rahtu. V-slowfast network for efficient visual sound separation. In *WACV*, 2022. 2

[93] Yizhe Zhu, Mohamed Elhoseiny, Bingchen Liu, Xi Peng, and Ahmed Elgammal. A generative adversarial approach for zero-shot learning from noisy texts. In *CVPR*, 2018. 2

[94] Yi Zhu, Yang Long, Yu Guan, Shawn Newsam, and Ling Shao. Towards universal representation for unseen action recognition. In *CVPR*, 2018. 2

[95] Yizhe Zhu, Jianwen Xie, Bingchen Liu, and Ahmed Elgammal. Learning feature-to-feature translator by alternating back-propagation for generative zero-shot learning. In *ICCV*, 2019. 2