



VGGSounder: Audio-Visual Evaluations for Foundation Models

Daniil Zverev^{1,*} Thaddäus Wiedemer^{2,3,4,*} Ameya Prabhu^{2,3}
 Matthias Bethge^{2,3} Wieland Brendel^{3,4} A. Sophia Koepke^{1,2,3}

¹Technical University of Munich, MCML ²University of Tübingen ³Tübingen AI Center
⁴MPI for Intelligent Systems, ELLIS Institute Tübingen

{daniil.zverev, a-sophia.koepke}@tum.de, {thaddaeus.wiedemer, ameya.prabhu}@uni-tuebingen.de

Abstract

The emergence of audio-visual foundation models underscores the importance of reliably assessing their multi-modal understanding. The VGGSound dataset is commonly used as a benchmark for evaluation audio-visual classification. However, our analysis identifies several limitations of VGGSound, including incomplete labelling, partially overlapping classes, and misaligned modalities. These lead to distorted evaluations of auditory and visual capabilities. To address these limitations, we introduce VGGSounder, a comprehensively re-annotated, multi-label test set that extends VGGSound and is specifically designed to evaluate audio-visual foundation models. VGGSounder features detailed modality annotations, enabling precise analyses of modality-specific performance. Furthermore, we reveal model limitations by analysing performance degradation when adding another input modality with our new modality confusion metric. Our dataset and project page are available at <https://vggsounder.github.io/>.

1. Introduction

Rigorous evaluation benchmarks have been instrumental in assessing the effectiveness of audio-visual models [33, 43, 49, 57]. Specifically, multi-modal foundation models integrating visual and auditory data aim to achieve a holistic understanding of audio-visual content. However, the field lacks large-scale modality-aware classification benchmarks with ground-truth annotations indicating whether each label is visible, audible, or both. Such annotations would allow detailed evaluations of multi-modal model capabilities. To address this gap, we introduce VGGSounder, an enhanced version of the widely-used audio-visual classification dataset VGGSound [13], which facilitates modality-

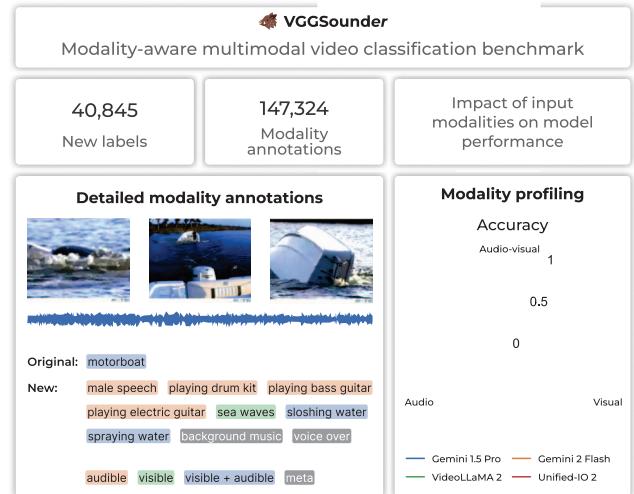


Figure 1. We introduce VGGSounder, a multi-label audio-visual classification benchmark with modality annotations. We extend the original VGGSound test set with human-annotated `audible`, `visible`, and `visible+audible` labels. We add `meta` labels for common confounders, such as background music. We benchmark eleven recent audio-visual models on VGGSounder. It enables selective analysis of a model’s auditory and visual capabilities on classes relevant for the queried modality.

aware evaluation of audio-visual foundation models.

VGGSound has several notable limitations. First, its data is inherently multi-label; for instance, a single sample may simultaneously include labels such as `playing drum kit` and `playing acoustic guitar` when multiple instruments are present. Additionally, evaluating how different modalities contribute to model performance becomes difficult without explicit modality annotations, as some labels are either not visually present or not audible (e.g., certain instruments might only be audible but not visible in advertisements). Moreover, overlapping label classes present another challenge; for example,

*equal contribution

the `orchestra` label often coincides with labels for individual instruments. These issues result in systematic under-evaluation of multi-modal audio-visual foundation models.

To overcome these limitations, we present VGGSounder, an improved benchmark inspired by similar advancements in other domains [12] [28]. We re-annotate the dataset to create a comprehensive multi-label classification setting by collecting detailed annotations for each sample, including (1) additional classes present, (2) explicit modality annotations to label modality misalignment, (3) metadata indicating the presence of background music, voice-over, or static images, and (4) merging of classes to address overlapping classes. Consequently, VGGSounder provides a robust, foundation-model-ready benchmark enabling structured analysis of whether models rely on audio or visual cues. Furthermore, we include meta-labels (e.g., background music, voice-over, or static images) to easily filter out unreliable labels during evaluation. Utilising VGGSounder, we evaluate audio-visual foundation models, demonstrating their poor performance on our benchmark. We find that the state-of-the-art, closed-source Gemini models consistently rely exclusively on the visual modality. In addition to that, we measure the modality confusion, i.e. when models get distracted by an additional input modality, which exposes the unsuccessful merging of modalities. These findings highlight the importance of the audio-focused VGGSounder benchmark as a critical tool for accurately assessing audio-visual foundation models.

We make the following contributions:

1. We illustrate limitations of VGGSound in Sec. 3
2. We curate VGGSounder with multi-modal human annotations for multi-label classification in Sec. 4
3. We evaluate state-of-the-art audio-visual models, observing differences between embedding models and autoregressive foundation models in Sec. 5
4. We propose new metrics to quantify the negative impact of using multiple input modalities in Sec. 5

2. Related work

Audio-visual learning Many prior works consider audio-visual tasks that include sound source localisation and separation [3, 5, 9, 15, 27, 56, 62, 67, 75, 80, 85, 86, 90], event localisation [50, 51, 74, 78], audio-visual question answering [48, 54, 83, 84], audio-visual synchronisation [14, 23, 25, 38, 39, 42], audio synthesis using visual information [19, 26, 31, 44, 45, 61, 69, 71, 87], or audio-driven face image synthesis [7, 40, 77]. Audio-visual data has also been leveraged for speech-related tasks, including speech and speaker recognition [2, 4, 59], or the spotting of spoken keywords [58, 66].

Furthermore, the natural alignment between audio and

Dataset	# Clips	# Classes	Multi-label	Modality labels	Annotation pipeline
Flickr-SoundNet [10]	2M	-	✗	✗	-
Kinetics-Sound [7]	18.8K	34	✗	✗	MTurk
AudioSet [26]	2.1M	537	✓	✗	Manual
└ AVE [62]	4K	28	✓	✗	Manual
└ VEGAS [76]	132K	10	✗	✗	MTurk
└ Visually Aligned Sounds [15]	13K	8	✗	✗	MTurk
VGGSound [12]	200K	309	✗	✗	Classifiers+Manual
└ VGGSound-Sparse [32]	7.1K	12	✗	✗	Manual
└ Visual Sound [66]	91K	309	✗	✗	ImageBind [30]
VGGSounder	15.4K	309	✓	✓	MTurk

Table 1. Comparison of audio-visual classification benchmarks.

video has been exploited to learn improved audio-visual embeddings for downstream tasks [6, 10, 11, 18, 20, 21, 46, 60, 63–65, 79]. Using both modalities jointly generally leads to performance boosts over using one modality in isolation. We examine this observation closely and aim to evaluate the effective use of multiple input modalities for the video classification task. To enable this, we propose — to the best of our knowledge, the first multi-label video classification benchmark that includes per-modality annotations for every sample (see Tab. 1).

Audio-visual foundation models Recently, multi-modal general-purpose models have emerged that can handle diverse downstream tasks without task-specific finetuning — also referred to as multi-modal foundation models. For instance, images or language were used as the bridge between modalities including audio, image, and text [30, 89]. Building on this, PandaGPT [72] leverages Vicuna [22] and ImageBind’s embedding space to train a general multi-modal model exclusively on image-text pairs. Unified-IO 2 [53] employs universal tokenisation to process audio, video, and text. VideoLLaMA2 [21] uses a Spatial-Temporal Convolution connector in the visual branch before projecting audio and visual information into the LLM input space. The recently introduced Ola model [52] advances omni-modal processing through progressive modality alignment, using video to bridge audio and visual information. The Gemini models [73] are closed-source multi-modal models that achieve impressive performance on diverse downstream tasks. We use VGGSounder to benchmark the audio-visual capabilities of the aforementioned models.

Audio-visual classification benchmarks Audio-visual classification is distinct from general video classification (e.g. on YouTube-8M [1]), as classes typically cover both audible and visible actions or events. Commonly used datasets for audio-visual classification include Kinetics-Sound [8] sourced from the Kinetics dataset [41], Flickr-SoundNet [11] scraped from Flickr, and AudioSet [29] and VGGSound [13], both sourced from YouTube. Kinetics-Sound features manual labels of human actions, but covers only 34 classes. Flickr-SoundNet is much larger, but only a small subset is labelled. Similarly, only a small fraction of the roughly 2M AudioSet samples are annotated and have aligned audio and video.

Figure 2. **Limitations of VGGSound.** We show video frames from videos in the VGGSound test set along with their annotated label (grey) to demonstrate various limitations. **A.** VGGSound samples are labelled with a single class, yet many videos contain multiple distinct classes. **B.** Additionally, many classes partially overlap or are ambiguous. **C.** Some samples are labelled with classes that are not present in one of the modalities (i.e., the labelled class is not visible or audible).

In contrast, VGGSound ensures audio-visual correspondence for around 200 000 samples and was curated with an automatic pipeline involving class-list generation, and auditory and visual content verification. The visual verification step ensures that a class is represented in the centre frame. The VEGAS dataset [88] provides better quality assurances for a small subset of AudioSet with only 10 classes. Visually Aligned Sounds [16] subsamples VEGAS and AudioSet after human verification, and Visual Sound [76] subsamples VGGSound using a multi-modal embedding model, both aiming for high audio-visual correspondence. Similarly, VGGSound-Sparse [37] is a subset of VGGSound with a focus on temporally and spatially sparse synchronisation signals (e.g., short loud noises). Overall, VGGSound strikes the best balance between size, generality, and annotations, making it a common benchmark for audio-visual classification. We update VGGSound to sustain its usability for the development of the next generation of multi-modal foundation models.

3. Limitations of VGGSound

Since we are interested in the VGGSound dataset for benchmarking audio-visual multi-modal models, our analysis focuses on the VGGSound test set,¹ which consists of 15 446 video clips, each 10s long and labelled with one of 309 classes. We qualitatively identify several limitations of the VGGSound annotations outlined below and in Fig. 2

Co-occurring classes While VGGSound’s visual verification aimed to minimise multiple classes co-occurring in a clip, we find that most samples nevertheless clearly contain multiple classes, see Fig. 2A. In some cases, classes are temporally separated, e.g., showing male speech, man speaking and then cutting to

¹Although these issues most likely also apply to the training set.

footage of firing cannon. Most often, classes co-occur at the same time, sometimes for the entire duration of the video clip. Overlapping classes are often related, such as different instruments in a band or orchestra, but can also be entirely unrelated, e.g., donkey, ass braying co-occurring with playing violin. As additional empirical evidence, we provide a co-occurrence matrix computed using CAV-MAE [33], a state-of-the-art audio-visual model, in Appendix D.

Overlapping classes The issue of co-occurring classes is exacerbated by many of the 309 automatically generated classes partially overlapping in their definition, as illustrated in Fig. 2B. We found two pairs of synonymous classes: timpani and tympani and dog barking and dog bow-wow. Additionally, some classes are strict subclasses of others, such as the gender-specific versions of cattle mooing: cow lowing and bull bellowing; or the more specific variants of people eating: people eating noodle and people eating apple. Finally, several classes commonly appear together, such as playing snare drum which is often played as part of a drum kit, or semantically similar concepts: running electric fan and air conditioning noise, and sloshing water and splashing water.

Modality misalignment Despite VGGSound’s auditory and visual content verification, we find that many of the annotated classes are not visible or not audible, as shown in Fig. 2C. A large fraction of videos contains background music, voice-over and narration, or other background sounds like bird chirping, tweeting or cricket chirping without a visible source. Similarly, some videos contain visible but inaudible cues for classes like sea waves. Static images and slide shows accompa-

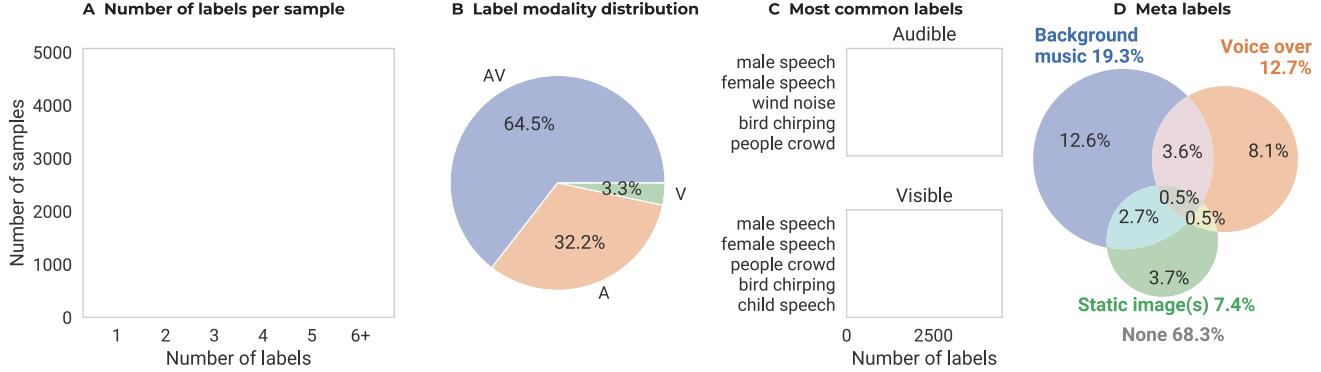


Figure 3. **Overview of VGGSounder.** **A.** Most samples contain more than one label. **B.** More than a quarter of labels are audible but not visible. In contrast, only a tiny fraction is visible but not audible. **C.** Speech and bird sounds are the most common classes; more details can be found in Appendix B. **D.** Forty percent of the samples contain some combination of background music, voice over, and static image(s), making the classification task significantly harder.

nied by music or other sounds are other frequent sources of misaligned modalities. Finally, some classes are misaligned by definition: `wind noise` is only audible and not visible. Overall, 48.43 % of the original VGGSound test samples have misaligned modalities. This finding challenges the widely held assumption that VGGSound has strong modality alignment [32, 47].

Other datasets, such as Visually Aligned Sounds, Visual Sound, and VGGSound-Sparse (see Sec. 2) omitted samples with misaligned modalities. In contrast, we contend that inaudible or invisible cues are common in natural videos and should be considered when benchmarking multi-modal models. We, therefore, place particular emphasis on crafting reliable modality annotations for all samples, allowing users to evaluate models on samples that guarantee modality alignment, and on those where classes are only visible or only audible (see Tab. I).

Takeaway 1 VGGSound suffers from several issues: class co-occurrence not captured by single labels, overlapping class definitions, and modality misalignment, see Fig. 2

4. Building VGGSounder

We propose a series of fixes for VGGSound’s issues, ultimately resulting in the updated VGGSounder benchmark. We are not the first aiming to future-proof an existing benchmark: [12] analysed shortcomings of ImageNet [24], ultimately proposing switching to a multi-label classification task with additional manual labels. [28] similarly re-annotated samples in MMLU [35, 36] to fix labelling errors. Both works inspire our approach to improve VGGSound. To deal with co-occurring classes, we switch to a multi-label classification setting. This effectively handles most overlapping class definitions: a strong model can assign a

high probability to multiple classes, even if they partially overlap. This also allows us to ensure that synonymous classes, as well as subclasses and their superclass, always appear together in the ground-truth labels.

To deal with modality misalignment, we add a modality annotation to each label. For example, we can label a video as containing `people clapping` in the audio and containing `playing volleyball` in the audio and video data. We also add meta-labels to indicate whether a sample contains `background music`, `voice over`, or `static image(s)` to optionally treat these cases separately during evaluation.

We employ a pipeline similar to [12] to annotate multiple labels per sample, which we outline below.

Collecting proposals We create a *gold standard* reference set by labelling a small, randomly selected subset of VGGSound test samples with four in-house computer vision experts. The interface used for this first labelling step is shown in Appendix A. We extend the subset until each class is covered at least once, leading to a final size of 417 samples. Labels from different annotators are merged via a simple majority vote.

Given the gold standard set, we want to find a solid strategy for automatically generating label proposals which are shown to the humans labelling the test set. This should have a recall greater than 90 % while maximising precision compared to the gold standard labels to produce a small set of proposals with good label coverage. Our final strategy combines predictions from several state-of-the-art models with a manual heuristic to obtain 93% recall for an average of 30 proposals per sample, see Appendix A.

Human labelling We use Amazon Mechanical Turk to re-annotate the entire VGGSound test set. For each sample, we first ask annotators to indicate whether the video contains `background music`, `voice over`, or

static image(s). Then, annotators are asked to indicate for each label proposal whether the class is `audible` and/or `visible`. Finally, annotators can add a class if it is missing from the proposals. Annotators were paid the US minimum wage; the interface used is shown in Appendix A. We let annotators label the samples in batches of 20, each containing two gold standard samples as catch trials. We reject and re-annotate all batches with a catch trial F_1 -score below 25%. In addition, we obtain modality annotations for the original VGGSound labels and meta-classes. Further details are provided in Appendix A.

Final labels We merge all obtained annotations using majority voting. Additionally, we automatically add synonymous classes and superclasses for a given subclass, e.g., we add `cattle mooing` whenever `cow lowing` is in the set of labels. The full set of classes added this way is described in Appendix A.

Takeaway 2 We develop VGGSounder: A multi-label video classification benchmark extending VGGSound with human-annotated multi-labels, modality annotations, and meta-labels as summarised in Fig. 3.

5. Benchmarking audio-visual models

We use VGGSounder to benchmark four popular audio-visual embedding models and seven foundation models, and analyse their auditory and visual capabilities.

Models We evaluate the audio-visual *embedding models* CAV-MAE [33], DeepAVFusion [57], AV-Siam [49], and Equi-AV [43]. Those were finetuned on VGGSound.

We benchmark several models from the closed-source Gemini family [73] in a zero-shot evaluation protocol. Furthermore, we use LLM-assisted evaluation to evaluate the following four open-source autoregressive *foundation models*: VideoLLaMA-2 [21], Unified-IO-2 [53], Panda-GPT [72], and Ola [52]. All models are evaluated in three modes: using unimodal-audio, unimodal-visual, or multi-modal (audio and visual) inputs. Further details about models and their evaluation are provided in Appendix C.1.

Metrics To benchmark the models on VGGSounder, we use multi-label classification metrics. For embedding models, all metrics are computed for the top- k predictions, with $k \in \{1, 3, 5, 10\}$. In contrast, prompting foundation models yields an unordered set of class predictions of varying size, and we compute only a single metric using the entire set. As a result, metrics are not directly comparable between embedding and foundation models. To get a sense of their relative performance, we report metrics for embedding models for $k = 1$ in the main text, matching the median number of predictions per sample for the foundation models.

For open-source models such as VideoLLaMA-2, Ola, Unified-IO-2, and Panda-GPT, we employ LLM-assisted

evaluation [55, 81], in which the Qwen-3 model [82] is tasked to assess the correspondence between model outputs and target classes. Closed-source models from the Gemini family are evaluated by providing the full list of 309 classes as input. Further details on the evaluation procedures and exact prompts are provided in Appendix C.

Subset accuracy compares the predicted label set to the ground-truth label set and reports the fraction of samples for which they match. This is our strictest metric.

F_1 -score is the harmonic mean of precision and recall. It is strictly larger than the subset accuracy.

Hit reports the fraction of samples for which *any* of the predicted labels are part of the ground-truth label set. This is the most lenient metric which is strictly larger than the F_1 -score. We include this metric for ease of comparison to the “Real-Accuracy” used in [12].

All metrics are computed separately for each input modality (audio, video, and audio-visual) and label modality. We use lowercase symbols a , v , and av to indicate the input modality: audio-only, visual-only, or audio-visual inputs, respectively. For label modality, we use uppercase symbols A , V , and AV , referring to the subsets of the benchmark with audible, visible, and audio-visual labels. We further include $A \neg V$ (audible but not visible) and $V \neg A$ (visible but not audible) to analyze unaligned cues. For clarity, we define shorthand notations such as $a = a(A)$ to denote the model’s performance on the audible subset A using only audio input. Analogously, $v = v(V)$ and $av = av(AV)$ refer to video-only and audio-visual performance on their respective label subsets. Furthermore, we use micro-aggregation to balance the contribution from each class.

We additionally measure the negative impact of using multimodal inputs. In particular, μ is a new metric we propose to measure a model’s *modality confusion* (μ). We define it as

$$\mu_M = 100 \cdot \frac{\sum_{x \in M} \mathcal{I}[m(x)\text{-correct} \cap av(x)\text{-wrong}]}{N_{total}}, \quad (1)$$

where $M \in [A, V, A \cap V]$ and their associated modality inputs are $m \in [a, v]$, correct/wrong is determined as in the *Hit* score (with $k = 1$ for embedding models). N_{total} refers to the total number of samples. μ measures the fraction of samples a model correctly classified given an input modality but got wrong when using both modalities simultaneously. We additionally report $\mu_{A \cap V}$ as the percentage of samples a model could solve in *either modality* unimodally but could not solve multi-modally. In other words, the modality confusion μ captures how frequently a model is distracted by an additional input modality, which can indicate the unsuccessful merging of modalities.

Takeaway 3 We propose a new metric, *modality confusion* μ , that measures how frequently a model is distracted by an additional input modality; see Eq. (1).

Model	Subset Accuracy \uparrow			$F_1 \uparrow$					Hit \uparrow			$\mu \downarrow$		
	a	v	av	a	v	av	$a(A \neg V)$	$v(V \neg A)$	a	v	av	μ_A	μ_V	$\mu_{A \cap V}$
<i>Embedding Models</i>														
CAV-MAE	13.19	19.23	24.49	34.46	34.91	42.62	13.94	19.00	62.29	53.44	64.17	3.58	6.43	0.77
DeepAVFusion	10.19	11.10	21.53	25.31	21.29	37.35	10.37	10.55	45.77	32.61	56.27	3.74	3.93	0.17
Equi-AV	11.60	10.52	20.00	29.39	20.42	34.69	12.55	10.65	53.12	31.26	52.24	6.97	7.13	1.38
AV-Siam	12.79	19.75	22.83	33.30	35.41	39.43	12.90	18.21	60.19	54.20	59.36	9.36	8.80	3.58
<i>Closed-source Foundation Models</i>														
Gemini 1.5 Flash	1.78	14.44	16.44	14.49	36.98	42.52	15.61	21.61	32.73	47.36	59.10	10.22	4.25	0.77
Gemini 1.5 Pro	3.05	20.86	22.53	19.26	49.73	53.74	17.73	22.90	35.03	69.23	75.42	2.09	4.85	0.57
Gemini 2.0 Flash	1.85	12.54	12.69	11.80	34.08	36.45	6.19	18.90	18.51	43.83	47.72	2.39	5.43	1.00
<i>Open-source Foundation Models</i>														
VideoLLaMA 2	12.86	19.85	24.47	38.87	47.82	52.35	20.34	28.08	58.91	52.02	59.80	12.72	5.46	2.95
Unified-IO 2	11.94	11.56	25.61	35.31	27.92	48.89	21.38	16.53	54.39	31.05	65.11	8.70	5.16	1.79
PandaGPT	3.19	4.19	5.46	18.73	18.56	20.85	16.82	14.40	21.08	17.01	18.82	7.59	5.90	2.47
OLA	14.11	8.69	18.19	47.70	24.85	46.48	40.44	13.45	59.05	24.57	51.51	15.47	6.80	2.49

Table 2. **Audio-visual video classification results on VGGSounder.** We report multi-label classification metrics (subset accuracy, F_1 -score, Hit accuracy, modality confusion (μ) for audio- $a(A)$, visual - $v(V)$, audio-visual - $av(AV)$, audio-only - $a(A \neg V)$ and video-only - $v(V \neg A)$ inputs. The embedding models CAV-MAE, DeepAVFusion, and Equi-AV were finetuned on the VGGSound training set. We report metrics for $k = 1$ here and for other k in Appendix D. The closed sourced multi-modal foundation models Gemini and open-sourced models use a zero-shot evaluation protocol and LLM-assisted protocol respectively.

5.1. Re-evaluating the state of the art

We present the benchmark performance of state-of-the-art audio-visual models in Tab. 2.

Overall performance Unsurprisingly, all models perform best with access to both input modalities (AV). Across all metrics, both open- and closed-source general-purpose models perform comparably to the purpose-built embedding models CAV-MAE and AV-Siam. This indicates that foundation models have reached — and for some modalities exceeded — the performance of specialised models. However, the embedding models finetuned on VGGSound generally have stronger unimodal performance with audio inputs (A) compared to visual inputs (V), which is in line with their pretraining on AudioSet. Interestingly, this trend is reversed for most foundation models, which seem to be biased towards visual inputs, with unimodal video performance (V) being substantially higher than unimodal audio performance (A).

Takeaway 4 Foundation models perform comparably to finetuned embedding models. Embedding models more heavily rely on audio cues than on visual ones, while foundation models exploit visible cues rather than audible ones, see Tab. 2.

Modality confusion The modality confusion score μ shows that, for all models, a notable fraction of test samples (4–11%) were misclassified when an additional modality was included—despite being correctly classified with unimodal input. Furthermore, for all models, a small portion of

test samples is not solvable multi-modally even though they were solvable in both modalities alone ($\mu_{A \cap V}$). This insight is made possible by VGGSounder’s per-label modality annotations and shows that all models are susceptible to being distracted given an additional modality. This is a concerning issue for multi-modal models since they should preserve unimodal capabilities when adding a second modality. Being able to evaluate this behaviour on the VGGSounder benchmark is a first step towards enabling the development of mitigation strategies, eventually resulting in stronger audio-visual models.

Takeaway 5 Our modality confusion score reveals that all models are negatively impacted by additional modalities for a substantial amount of samples (see Tab. 2).

Performance across modalities Fig. 4 shows the performance profiles across modalities. At first glance, we can see that VideoLLaMA-2’s performance is well balanced for different input modalities, while models from the Gemini family distinctly underperform on audio inputs. In contrast, embedding models exhibit a moderate balance across modalities, with DeepAVFusion and EquiAV showing slight underperformance for visual input.

As Fig. 4 also illustrates, profiling of this kind is enabled through the modality annotations in VGGSounder. In contrast, VGGSound assumed that all classes are perceptible in both modalities, and did not account for background sounds. This resulted in consistent under-evaluation of foundation models (that were not finetuned on VGGSound) for audio inputs.

In addition to the radar plot in Fig. 4, we provide results on VGGSound in Appendix D showing that all models have substantially lower performance than their hit scores in Tab. 2. This confirms that many model predictions were incorrectly flagged as false positives in VGGSound due to the incomplete ground-truth labels, painting a distorted picture of models’ limitations.

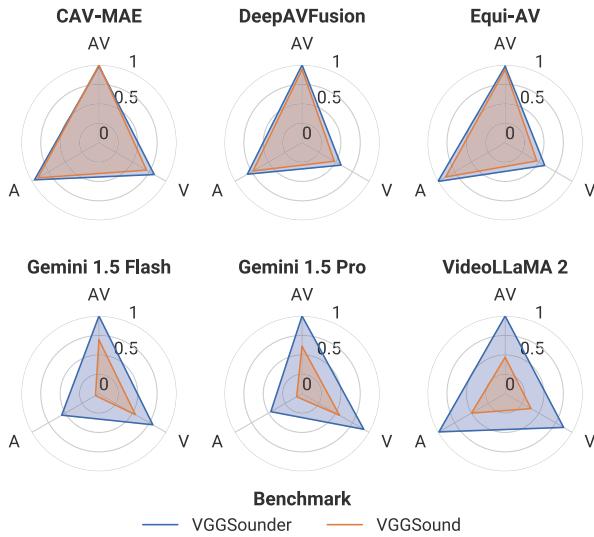


Figure 4. VGGSounder more accurately captures model performance across input modalities. We show the Hit score on VGGSounder and accuracy on VGGSound, normalised by the per-model maximum performance on each benchmark. Specifically for foundation models, we observe a significant difference in performance between VGGSound and VGGSounder.

Takeaway 6 VGGSounder’s more complete ground-truth labels allow for more accurate, modality-specific profiling of model performance (see Fig. 4).

5.2. Performance analysis using meta-classes

VGGSounder includes annotations of three meta-classes for each sample: `background music` indicates whether samples contain music without a visible source, `voice over` similarly marks speech without a visible source, and `static image(s)` flags that the visual stream consists of one or only a few static visual frames.

These new meta-classes allow us to evaluate the model behaviour in challenging scenarios where information from one modality dominates. We consider the performance difference between samples that *do* contain a meta-label and samples that *do not* contain it. Positive numbers indicate that the models perform better on the subset with meta-labels. In Tab. 3, we summarise the main findings, focussing on F_1 -score as the most balanced multi-label metric. Additional results are provided in Appendix D.

Background music All models perform worse on video samples containing background music. This indicates that it is challenging to decouple background audio from the rest of the video. The evaluated models are not good at differentiating between sound sources without visual cues, e.g. predicting different instruments in the background music.

Voice over In contrast to background music, we observe a clear difference between embedding models and foundation models for samples with voice-over. While the audio classification performance of embedding models drops substantially. This drop is only slight for VideoLLaMA-2 and Unified-IO 2, and the performance of other foundation models even improves. This indicates that the foundation models are less distracted by voice-over.

Static image(s) The impact of static images is more nuanced: First, audio classification performance improves for embedding models while it decreases for the foundation models. This shows that the purpose-built, VGGSound-finetuned embedding models can more accurately predict specific sounds in the absence of other cues.

Second, visual classification performance on static images drops for all models, suggesting that models rely on rich temporal cues to make accurate predictions.

Third, the subset accuracy (i.e. exact label set matches) for uni-modal predictions shows that foundation models perform better in visual than audio classification with and without static images. In contrast, embedding models perform better in audio classification on static images but worse in visual classification for other samples. This suggests that *non-static* samples form a challenging subset for embedding models, where a model needs to favour one modality over another one to make a correct prediction.

Samples without any meta-label When comparing the model performance on samples without background music, static images, or voice-over annotations (column *neither* in Tab. 3) to the performance on the entire test set, we see a performance gain (here negative numbers indicate an increased F_1 score). This finding concludes that these three categories form challenging subsets of the dataset.

Takeaway 7 Samples with `background music`, `static image(s)`, and `voice over` provide distinct challenges for each model (see Tab. 3). This highlights VGGSounder’s value for comprehensive model evaluation.

5.3. Impact of VGGSounder labels

Our relabelling pipeline adds two types of labels to those in VGGSound: (1) automatically generated labels based on synonymous classes and sub-/superclasses, and (2) human-curated labels. In Tab. 4, we ablate the impact of each type of added labels in terms of the relative performance gains (Hit score). A complete breakdown of the effects across all

Model	background music			voice over			static image(s)						neither		
	ΔF_1			ΔF_1			ΔF_1		Sub. Acc. w/		Sub. Acc. w/o		ΔF_1		
	a	v	av	a	a	v	a	v	a	v	a	v	a	v	av
<i>Embedding Models</i>															
CAV-MAE	-0.66	-0.59	-0.65	-1.13	0.31	-0.28	12.73	19.22	11.98	19.21	-1.00	-1.27	-1.24		
DeepAVFusion	-0.71	-0.80	-0.70	-1.25	0.28	-0.18	9.79	10.81	9.30	10.81	-1.27	-1.01	-1.36		
Equi-AV	-0.79	-0.42	-0.41	-0.97	0.28	-0.23	11.12	10.45	10.49	10.45	-0.97	-0.97	-0.89		
AV-Siam	-0.73	-0.72	-0.82	-1.08	0.34	-0.26	12.34	19.53	11.57	19.52	-1.01	-1.34	-1.39		
<i>Closed-source Foundation Models</i>															
Gemini 1.5 Flash	-0.22	-0.41	-0.69	2.08	-0.36	-0.25	1.66	14.35	1.68	14.39	1.40	-1.10	-1.62		
Gemini 1.5 Pro	-0.36	-0.61	-0.94	2.48	-0.31	-0.24	2.83	20.85	2.87	20.80	1.64	-0.89	-1.09		
Gemini 2.0 Flash	-0.09	-0.32	-0.58	0.03	0.13	-0.26	1.70	12.33	1.53	12.40	-0.03	-1.07	-1.22		
<i>Open-source Foundation Models</i>															
VideoLLaMA 2	-0.47	-0.72	-0.84	-0.54	0.27	-0.31	12.55	19.64	12.04	19.70	-0.47	-1.26	-1.33		
Unified-IO 2	-1.20	0.18	-0.65	-0.68	0.12	-0.21	11.39	11.89	10.88	12.00	-1.07	-0.15	-1.17		
PandaGPT	-1.13	-0.15	-0.42	0.56	-0.22	-0.16	2.94	4.27	2.91	4.24	-0.60	-0.15	-0.65		
OLA	-2.16	0.09	-0.44	1.50	-0.48	-0.18	13.03	8.88	12.88	8.89	-0.80	-0.13	-0.06		

Table 3. **Summary of performance differences in the presence/absence of meta-classes.** Difference in F_1 scores (ΔF_1) for audio-visual video classification on VGGSounder between videos with a meta-class and those without it. Positive numbers (Δ) indicate better performance when the meta-class is present. Additional results are provided in Appendix D

Model	Human labels \uparrow			Auto labels \uparrow		
	A	V	AV	A	V	AV
Gemini 1.5 Flash	29.28	14.59	16.36	0.48	0.93	1.51
Gemini 1.5 Pro	28.61	25.52	27.63	0.31	1.99	2.10
Gemini 2.0 Flash	8.80	12.16	11.13	0.22	1.12	1.28

Table 4. **Impact of added labels using different strategies in VGGSounder.** We show the change in multi-label classification performance (Δ Hit) when adding automatically added (Auto) or human-annotated (Human) labels to VGGSound, and compare to the original VGGSound data.

models and metrics is provided in Appendix D. While performance is consistently higher with automatically added labels (Auto), the increase is noticeably smaller than that for human-curated labels. Paired with the observation that models do frequently predict correct classes that were not part of the original VGGSound label set, this indicates that human-curated labels better cover the ground truth.

Takeaway 8 Automatically added labels are an important step, but human-curated labels have a bigger effect on eliminating incorrectly flagged false positives, underscoring the value of accurate human annotation.

6. Discussion

Choice of VGGSound as base dataset VGGSound is commonly used to evaluate audio-visual models on the multi-modal video classification task. As it is currently the most suitable testbed for audio-visual classification tasks (due to its size, diversity of categories, non-constrained setting, and relatively strong audio-visual correspondence), it serves as an optimal starting point for our substantially improved VGGSounder benchmark with a multi-label evaluation pro-

tocol for foundation models that makes the benchmark suitable for meaningful evaluation.

Multi-label vs single-label classification Video content is inherently complex, often containing multiple co-occurring objects and actions both within and across modalities. This makes it unlikely that a given clip belongs to just one class as is the case in the single-label classification task. Therefore, our VGGSounder extends the VGGSound test set to multi-label classification. Unlike models trained on a narrow single-label dataset, foundation models develop versatile representations. Multi-label evaluation thus also aligns very naturally with this capability.

7. Conclusion

We introduced an modality-aware evaluation set for audio-visual foundation models. VGGSounder builds on the widely used VGGSound dataset by adding: (a) comprehensive human annotations for missing classes, (b) specifying modality information per label, (c) introducing specialised meta-labels for frequently occurring real-world challenges, and (d) using heuristic methods to improve label quality. Through our newly introduced metric, modality confusion, we observe that incorporating additional modalities does not necessarily yield better results. Models often become more confused on a substantial subset of test samples. Furthermore, finetuned embedding models tend to rely heavily on audio cues, while foundation models depend more on visual information. Additionally, our meta-label analysis highlights distinct challenges across various specialised yet commonly occurring scenarios such as background music, static images, and voice-overs. Overall, we hope that the VGGSounder benchmark will advance the evaluation and development of foundational audio-visual models.

Acknowledgements

The authors would like to thank Felix Förster, Sayak Mallick, and Prasanna Mayilvahanan for their help with data annotation, Felix Förster and Monica Riedler for proofreading the paper, and Thomas Klein and Shyamgopal Karthik for their help in setting up MTurk. They also thank numerous MTurk workers for labelling.

This work was in part supported by the BMBF (FKZ: 01IS24060, 01I524085B, 01IS18039A), the DFG (SFB 1233, project number: 276693517), and the Open Philanthropy Foundation funded by the Good Ventures Foundation. WB acknowledges financial support via an Emmy Noether Grant funded by the German Research Foundation (DFG) under grant no. BR 6382/1-1. WB is a member of the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645. This research utilised compute resources at the Tübingen Machine Learning Cloud, DFG FKZ INST 37/1057-1 FUGG. The authors thank the IMPRS-IS for supporting TW.

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.
- [2] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. Deep audio-visual speech recognition. *IEEE TPAMI*, 2018.
- [3] Triantafyllos Afouras, Yuki M Asano, Francois Fagan, Andrea Vedaldi, and Florian Metze. Self-supervised object detection from audio-visual correspondence. In *ECCV*, 2020.
- [4] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. Asr is all you need: Cross-modal distillation for lip reading. In *ICASSP*, 2020.
- [5] Triantafyllos Afouras, Andrew Owens, Joon Son Chung, and Andrew Zisserman. Self-supervised learning of audio-visual objects from video. In *ECCV*, 2020.
- [6] Humam Alwassel, Dhruv Mahajan, Lorenzo Torresani, Bernard Ghamem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. In *NeurIPS*, 2020.
- [7] Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. Facetalk: Audio-driven motion diffusion for neural parametric head models. In *CVPR*, 2024.
- [8] Relja Arandjelović and Andrew Zisserman. Look, listen and learn. In *ICCV*, 2017.
- [9] Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *ECCV*, 2018.
- [10] Yuki M. Asano, Mandela Patrick, Christian Rupprecht, and Andrea Vedaldi. Labelling unlabelled videos from scratch with multi-modal self-supervision. In *NeurIPS*, 2020.
- [11] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *NeurIPS*, 2016.
- [12] Lucas Beyer, Olivier J Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020.
- [13] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggssound: A large-scale audio-visual dataset. In *ICASSP*, 2020.
- [14] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Audio-visual synchronisation in the wild. In *BMVC*, 2021.
- [15] Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *CVPR*, 2021.
- [16] Peihao Chen, Yang Zhang, Mingkui Tan, Hongdong Xiao, Deng Huang, and Chuang Gan. Generating visually aligned sound from videos. *IEEE Transactions on Image Processing*, 2020.
- [17] Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*, 2022.
- [18] Yanbei Chen, Yongqin Xian, A. Sophia Koepke, Ying Shan, and Zeynep Akata. Distilling audio-visual knowledge by compositional contrastive learning. In *CVPR*, 2021.
- [19] Ziyang Chen, Prem Seetharaman, Bryan Russell, Oriol Nieto, David Bourgin, Andrew Owens, and Justin Salamon. Video-guided foley sound generation with multimodal controls. In *CVPR*, 2025.
- [20] 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.
- [21] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-lmms. *arXiv preprint arXiv:2406.07476*, 2024.
- [22] Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. 2023.
- [23] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *ACCV*, 2016.
- [24] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [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.
- [26] Chuang Gan, Deng Huang, Peihao Chen, Joshua B Tenenbaum, and Antonio Torralba. Foley music: Learning to generate music from videos. In *ECCV*, 2020.
- [27] Ruohan Gao and Kristen Grauman. Co-separating sounds of visual objects. In *ICCV*, 2019.
- [28] Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad

- Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kaddour, Emile van Krieken, and Pasquale Minervini. Are we done with mmlu? *arXiv preprint arXiv:2406.04127*, 2024.
- [29] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *ICASSP*, 2017.
- [30] Rohit Girdhar, Alaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023.
- [31] Shir Goldstein and Yael Moses. Guitar music transcription from silent video. In *BMVC*, 2018.
- [32] Yuan Gong, Alexander H Liu, Andrew Rouditchenko, and James Glass. Uavm: Towards unifying audio and visual models. *IEEE Signal Processing Letters*, 2022.
- [33] Yuan Gong, Andrew Rouditchenko, Alexander H Liu, David Harwath, Leonid Karlinsky, Hilde Kuehne, and James Glass. Contrastive audio-visual masked autoencoder. In *ICLR*, 2023.
- [34] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.
- [35] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. In *ICLR*, 2021.
- [36] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *ICLR*, 2021.
- [37] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Sparse in space and time: Audio-visual synchronisation with trainable selectors. In *BMVC*, 2022.
- [38] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Sparse in space and time: Audio-visual synchronisation with trainable selectors. In *BMVC*, 2022.
- [39] Vladimir Iashin, Weidi Xie, Esa Rahtu, and Andrew Zisserman. Synchformer: Efficient synchronization from sparse cues. In *ICASSP*, 2024.
- [40] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman. You said that?: Synthesising talking faces from audio. *IJCV*, 2019.
- [41] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [42] Naji Khosravan, Shervin Ardeshir, and Rohit Puri. On attention modules for audio-visual synchronization. In *CVPR Workshop*, 2019.
- [43] Jongsuk Kim, Hyeongkeun Lee, Kyeongha Rho, Junmo Kim, and Joon Son Chung. Equiav: Leveraging equivariance for audio-visual contrastive learning. In *ICML*, 2024.
- [44] A. Sophia Koepke, Olivia Wiles, and Andrew Zisserman. Visual pitch estimation. In *SMC*, 2019.
- [45] 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.
- [46] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, 2018.
- [47] Sangho Lee, Jiwan Chung, Youngjae Yu, Gunhee Kim, Thomas Breuel, Gal Chechik, and Yale Song. Acav100m: Automatic curation of large-scale datasets for audio-visual video representation learning. In *ICCV*, 2021.
- [48] Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer questions in dynamic audio-visual scenarios. In *CVPR*, 2022.
- [49] Yan-Bo Lin and Gedas Bertasius. Siamese vision transformers are scalable audio-visual learners. In *ECCV*, 2024.
- [50] Yan-Bo Lin and Yu-Chiang Frank Wang. Audiovisual transformer with instance attention for audio-visual event localization. In *ACCV*, 2020.
- [51] Yan-Bo Lin, Yu-Jhe Li, and Yu-Chiang Frank Wang. Dual-modality seq2seq network for audio-visual event localization. In *ICASSP*, 2019.
- [52] Zuyan Liu, Yuhao Dong, Jiahui Wang, Ziwei Liu, Winston Hu, Jiwen Lu, and Yongming Rao. Ola: Pushing the frontiers of omni-modal language model. *arXiv preprint arXiv:2502.04328*, 2025.
- [53] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *CVPR*, 2024.
- [54] Jie Ma, Min Hu, Pinghui Wang, Wangchun Sun, Lingyun Song, Hongbin Pei, Jun Liu, and Youtian Du. Look, listen, and answer: Overcoming biases for audio-visual question answering. *NeurIPS*, 2024.
- [55] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- [56] Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In *ECCV*, 2022.
- [57] Shentong Mo and Pedro Morgado. Unveiling the power of audio-visual early fusion transformers with dense interactions through masked modeling. In *CVPR*, 2024.
- [58] Liliane Momeni, Triantafyllos Afouras, Themis Stafylakis, Samuel Albanie, and Andrew Zisserman. Seeing wake words: Audio-visual keyword spotting. In *BMVC*, 2020.
- [59] Arsha Nagrani, Joon Son Chung, Samuel Albanie, and Andrew Zisserman. Disentangled speech embeddings using cross-modal self-supervision. In *ICASSP*, 2020.
- [60] Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. In *NeurIPS*, 2021.
- [61] 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.

- [62] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, 2018.
- [63] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *ECCV*, 2016.
- [64] 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.
- [65] 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.
- [66] KR Prajwal, Liliane Momeni, Triantafyllos Afouras, and Andrew Zisserman. Visual keyword spotting with attention. In *BMVC*, 2021.
- [67] Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyo Lin. Multiple sound sources localization from coarse to fine. In *ECCV*, 2020.
- [68] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavy, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *ICLR*, 2023.
- [69] Kun Su, Xiulong Liu, and Eli Shlizerman. Multi-instrumentalist net: Unsupervised generation of music from body movements. *arXiv preprint arXiv:2012.03478*, 2020.
- [70] Kun Su, Xiulong Liu, and Eli Shlizerman. How does it sound? In *NeurIPS*, 2021.
- [71] Kun Su, Xiulong Liu, and Eli Shlizerman. From vision to audio and beyond: A unified model for audio-visual representation and generation. In *ICML*, 2024.
- [72] Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- [73] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [74] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chen-liang Xu. Audio-visual event localization in unconstrained videos. In *ECCV*, 2018.
- [75] 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.
- [76] Ilpo Viertola, Vladimir Iashin, and Esa Rahtu. Temporally aligned audio for video with autoregression. In *ICASSP*, 2025.
- [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.
- [78] Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention matching for audio-visual event localization. In *CVPR*, 2019.
- [79] Fanyi Xiao, Yong Jae Lee, Kristen Grauman, Jitendra Malik, and Christoph Feichtenhofer. Audiovisual slowfast networks for video recognition. *arXiv preprint arXiv:2001.08740*, 2020.
- [80] 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.
- [81] Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointlm: Empowering large language models to understand point clouds. In *ECCV*, 2024.
- [82] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [83] Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. Avqa: A dataset for audio-visual question answering on videos. In *ACM MM*, 2022.
- [84] Heeseung Yun, Youngjae Yu, Wonsuk Yang, Kangil Lee, and Gunhee Kim. Pano-avqa: Grounded audio-visual question answering on 360deg videos. In *ICCV*, 2021.
- [85] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018.
- [86] Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In *ICCV*, 2019.
- [87] Hang Zhou, Ziwei Liu, Xudong Xu, Ping Luo, and Xiaogang Wang. Vision-infused deep audio inpainting. In *ICCV*, 2019.
- [88] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L. Berg. Visual to sound: Generating natural sound for videos in the wild. In *CVPR*, 2018.
- [89] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *arXiv preprint arXiv:2310.01852*, 2023.
- [90] Lingyu Zhu and Esa Rahtu. V-slowfast network for efficient visual sound separation. In *WACV*, 2022.