

AVSET-10M: AN OPEN LARGE-SCALE AUDIO-VISUAL DATASET WITH HIGH CORRESPONDENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent research initiatives such as ChatGPT and Sora highlight the important role of large-scale data in advancing generative and comprehension tasks. However, the scarcity of comprehensive and large-scale audio-visual correspondence datasets poses a significant challenge to research in the audio-visual field. To address this gap, we introduce **AVSET-10M**, a high-correspondence audio-visual dataset comprising 10 million samples, featuring the following key attributes: (1) **High Audio-Visual Correspondence**: Through meticulous sample filtering, we ensure a strong correspondence between the audio and visual components of each entry. (2) **Comprehensive Categories**: Encompassing 527 unique audio categories, AVSET-10M provides a wide range of audio categories for diverse research needs. (3) **Large Scale**: With 10 million samples, AVSET-10M is one of the largest publicly available audio-visual correspondence datasets. We have benchmarked two critical tasks on AVSET-10M: audio-video retrieval and vision-queried sound separation. These tasks underscore the importance of precise audio-visual correspondence in advancing audio-visual research. For more information, please visit our demo page at <https://avset-10m.github.io/>.

1 INTRODUCTION

Scaling up significantly enhances performance in understanding (Touvron et al., 2023; Bai et al., 2023; Liu et al., 2024) and generation (Kondratyuk et al., 2023; Kang et al., 2023; Xiang et al., 2024) tasks across visual and language modalities. Inspired by the success of ImageNet (Deng et al., 2009) in visual research, some introduce the pioneering large-scale audio dataset, AudioSet (Gemmeke et al., 2017), which comprises 2.1 million audio samples each manually annotated with fine-grained audio categories to advance automatic audio understanding. However, the annotation process in AudioSet primarily focuses on only audio labels, neglecting the audio-visual correspondence. To address the need for exploring temporal consistency between audio and video, researchers develop the VGGSound (Chen et al., 2020), which includes 200,000 samples with audio-visual correspondence. Leveraging this dataset, significant breakthroughs have been achieved in the audio-visual domain, including vision-queried sound separation (Dong et al., 2022) and vision-based audio synthesis (Huang et al., 2023; Xing et al., 2024).

Meanwhile, the scale of vision-language datasets (Thomee et al., 2016; Miech et al., 2019; Xue et al., 2022; Schuhmann et al., 2022; Wang et al., 2023) has expanded dramatically, encompassing up to 100 million or even 1 billion samples. This expansion has facilitated a qualitative leap in understanding (Touvron et al., 2023; Liu et al., 2024) and generation (Kondratyuk et al., 2023) tasks within the vision and language fields, enabling the development of intelligent large language models (Touvron et al., 2023) and video generation technologies (Brooks et al., 2024) that simulate real-world scenarios. In contrast, the scale of datasets that ensure audio-visual correspondence remains markedly limited, posing a constraint on advancements in audio-visual field.

To further expand the audio-visual corresponding dataset and promote research on audio-visual temporal consistency, we propose AVSET-10M, the first 10 million scale audio-visual corresponding dataset, along with AVSET-700K, a subset containing fine-grained audio annotations. In Table 1, we present a comparison among various existing audio and audio-visual datasets. Our dataset construction process includes four stages: (1) Data collection, (2) Audio-visual correspondence filtering, (3) Voice-over filtering, and (4) Sample recycling with sound separation. AudioSet (Gemmeke et al., 2017),

Table 1: Comparison of different audio-video datasets. **AV-C** denotes the audio-visual correspondence. **# Class**: Number of audio categories. ACAV-100M[†] does not filter out the voiceover.

Datasets	Video	AV-C	#Class	#Clips	#Dur.(hrs)	#Avg Dur.(s)
DCASE2017 (Mesaros et al., 2019)	✗	✗	17	57K	89	3.9
FSD (Fonseca et al., 2017)	✗	✗	398	24K	119	17.4
AudioSet (Gemmeke et al., 2017)	✓	✗	527	2.1M	5.8K	10
AudioScope-V2 (Tzinis et al., 2022)	✓	✗	-	4.9M	1.6K	5
ACAV100M(Lee et al., 2021) [†]	✓	✗	-	100M	277.7K	10
HD-VILA-100M (Xue et al., 2022)	✓	✗	-	103M	371.5K	13.4
Panda-70M (Chen et al., 2024)	✓	✗	-	70.8M	166.8K	8.5
AVE (Tian et al., 2018)	✓	✓	28	4K	11	10
VGGSound (Chen et al., 2020)	✓	✓	309	200K	550	10
AVSET-700K (ours)	✓	✓	527	728K	2.0K	10
AVSET-10M (ours)	✓	✓	527	10.9M	30.4K	10.3

known for its fine-grained manual labeling of audio categories, is selected as our initial data source and develop AVSET-700K with accurate audio labels. To increase the number of samples per audio category, we choose Panda-70M (Chen et al., 2024) as an additional data source, expanding AVSET-700K to 10 million audio-visual corresponding samples. Panda-70M processes long videos into multiple semantically coherent sub-segments, effectively preventing the mixing of sounds from different events. Previous filtering method (Chen et al., 2020) using visual classification models struggles to distinguish audio events that cannot be identified by unique visual content, such as silence, thereby limiting the diversity of audio categories. To address this issue, we introduce a new filtering method based on audio-visual similarity (Girdhar et al., 2023), which significantly broadens the diversity of audio types. We employ an audio classification model (Kong et al., 2020) to filter out samples containing narration or background music that does not align with the visual content. As speech is commonly found in wild video data, which often results in the inadvertent filtering out of a substantial amount of audio samples containing voice-overs. This leads to the loss of many potentially useful and valuable samples across various audio categories. Thus, we further attempt to employ sound separation models (Solovyev et al., 2023) to recycle as many of these wasted samples as possible. From the initial 41 million samples, we filter 10 million audio-visual samples with high correspondence. Verification experiments demonstrate that our AVSET-700K provides more robust audio-visual correspondence than the previously used audio-visual corresponding dataset (VGGSound). Additionally, benchmarks of audio-video retrieval and vision-queried sound separation on AVSET-10M demonstrate it offers more research opportunities in the field of audiovisual studies.

2 RELATED WORKS

2.1 AUDIO-VISUAL MODELS

As multi-modal research progresses, the investigation (Li et al., 2022; Rahman et al., 2019; Ibrahim et al., 2023) into the correlations between audio and visual modalities has advanced. Initially, researchers employ both audio and video data to provide semantically richer information, thereby improving video understanding and significantly enhancing performance in various video understanding tasks such as video question answering (VQA) (Li et al., 2022; Akbari et al., 2021), video captioning (Rahman et al., 2019; Iashin & Rahtu, 2020a;b; Lin et al., 2023), and video retrieval (Lew et al., 2006; Ibrahim et al., 2023; Arora et al., 2024). Following these developments, ImageBind (Girdhar et al., 2023) emerges as a pioneering project that successfully aligns audio and visual content, marking a significant step in exploring semantic alignment between these modalities. Building on this foundation, subsequent research has delved into more intricate interactions between audio and video, achieving milestones in vision-queried sound separation (Dong et al., 2022) and video dubbing (Huang et al., 2023). However, while these methods have managed to align audio and visual content semantically, they often falter in maintaining temporal consistency. Some of the recent innovations (Luo et al., 2024) have introduced audio-visual temporal consistency supervision loss to enhance the temporal alignment in video dubbing.

108 Despite these advancements, the limited availability of training data continues to pose a significant
109 challenge, keeping the development of audio-visual temporal consistency at a rudimentary level. As
110 a result, the understanding of visual content remains largely confined to the semantic level, which
111 hampers the ability of models to accurately capture the audio-visual temporal consistency.

112 2.2 AUDIO-VIDEO DATASET

113 Inspired by ImageNet (Deng et al., 2009), researchers (Gemmeke et al., 2017) annotate a substantial
114 audio dataset, consisting of 2.1 million audio samples, aimed at enhancing automatic audio compre-
115 hension. Although annotators are encouraged to consult video content to refine the accuracy of audio
116 annotations, the dataset primarily focuses on precise audio annotations without additional measures
117 to filter out audio-visual non-corresponding samples. This limits the exploration of audio-video
118 consistency.

119 To investigate audio-visual consistency, researchers (Chen et al., 2020) employed a visual model
120 to identify sound-producing objects in videos, leading to the creation of VGGSound, a dataset
121 comprising 200,000 audio-visual corresponding samples. However, this visual model is effective
122 only in scenes characterized by definite actions or visible objects. It struggles to handle audio
123 events that lack distinctive visual content, such as silence and ambient sounds in urban outdoor
124 environments, even though there is a significant correlation between these audio events and the visual
125 elements in these scenes (e.g., silence audio in aquariums video). This limitation constrains the
126 diversity of audio categories represented in VGGSound. To further scale up audio-visual datasets,
127 ACAV100M (Lee et al., 2021) employs a clustering-based approach for data filtering. However,
128 it does not filter out voice-overs, resulting in the audio-visual correspondence of the final dataset
129 being even worse than that of AudioSet. AudioScope V1/2 (Tzinis et al., 2020; 2022) utilizes an
130 unsupervised audio-video consistency prediction model to evaluate audio-video matching scores,
131 screening 2,500 hours of video samples from YFCC100M (Thomee et al., 2016). Nevertheless, due
132 to limitations in prediction accuracy, the consistency between audio and video cannot be guaranteed,
133 and there remains a significant amount of inconsistent audio-visual content in the dataset.

134 Although subsequent research introduces larger video datasets (Xue et al., 2022; Wang et al., 2023;
135 Chen et al., 2023; 2024), the primary focus remains on exploring the relationship between video and
136 text, overlooking the audio-visual correspondence. To the best of our knowledge, our AVSET-10M
137 represents the largest open audio-visual high-correspondence dataset currently available, contain-
138 ing 10 million data samples across 527 different audio categories. This dataset opens up more
139 opportunities for research in the audio-video field.

140 3 AVSET-10M

141 3.1 DATASET CONSTRUCTION PIPELINE

142 **Stage 1: Data Collection.** We select two different open-source datasets, AudioSet (Gemmeke et al.,
143 2017) and Panda-70M (Chen et al., 2024), as data sources. All videos are sourced from open-domain
144 YouTube content. Since these datasets do not focus on audio-visual correspondence, they contain
145 a substantial number of mismatched audio-visual samples. We propose a filtering process to select
146 samples with high audio-visual correspondence, thereby constructing AVSET-10M.

147 AudioSet (Gemmeke et al., 2017) is a pioneering large-scale audio dataset where all audio category
148 labels are carefully annotated by human annotators. During the annotation process, annotators are
149 allowed to view the accompanying videos, which aids in accurate audio category identification. This
150 dataset includes 2.1 million audio samples across 527 unique audio categories. From AudioSet,
151 we select 727,530 samples that demonstrate high audio-visual correspondence with reliable audio
152 category labels to form AVSET-700K.

153 Additionally, to further expand the number of samples in each audio class, we select Panda-70M (Chen
154 et al., 2024), a large-scale video-text dataset containing 70 million semantically consistent segments.
155 It employs shot boundary detection technology (pyd) to divide the original videos into smaller
156 semantically consistent segments. This segmentation ensures that each clip contains only a single
157 event, preventing sound category conversion due to event switching and facilitating the subsequent
158
159
160
161

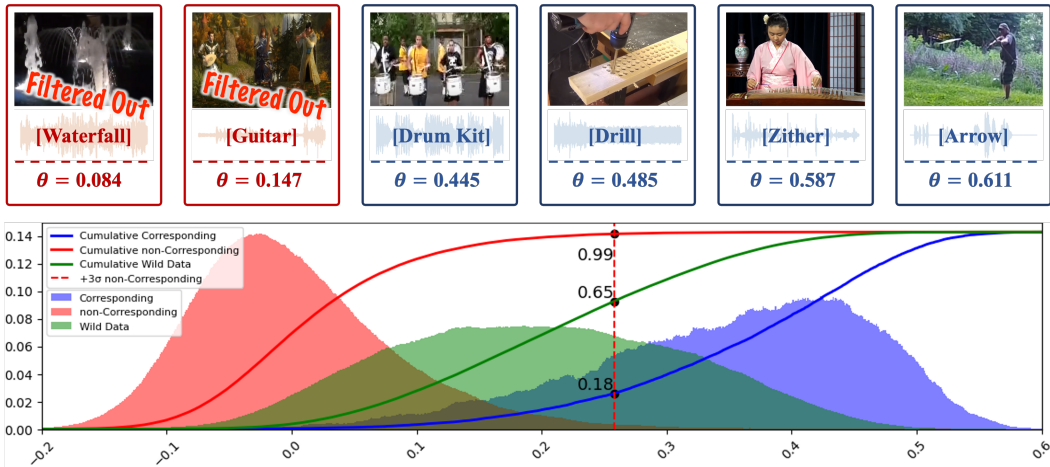


Figure 1: The distribution of audio-visual similarity among audio-visual corresponding samples, audio-visual non-corresponding samples and randomly selected wild samples. The similarity of non-corresponding data follows the distribution $N_{non-corresponding}(0.015, 0.081^2)$. Approximately 65% of the randomly selected wild samples and 18% of the audio-visual corresponding samples exhibit similarities below the $\mu + 3\sigma$ (0.2564) threshold of $N_{non-corresponding}$, suggesting a potential for these samples to be classified as audio-visual non-corresponding.

filtering process. Leveraging Panda-70M, we expand AVSET-700K to a total of 10 million audio-visual corresponding samples, thus forming AVSET-10M.

Stage 2: Audio-Visual Correspondence Filtering. Previous researchers (Chen et al., 2020) compute the cosine similarity between textual class label and visual content to gauge alignment confidence between vision and language. They subsequently filter video samples for each class label using a manually selected threshold. However, this method is effective only in scenes featuring definite actions or visual objects. It struggles to handle audio events that lack distinctive visual content, such as silence and urban outdoor environments, even though there is a significant correlation between these audio events and the visual elements in these scenes (e.g., video of aquariums and the silence audio). This consequently restricts the diversity of audio categories available in the dataset. We propose determining the confidence of audio-visual correspondence based on audio-visual similarity. This approach enables the screening of audio samples that lack distinctive visual content, thereby enhancing the diversity of samples in the dataset. Specifically, we randomly select 7,500 audio-visual corresponding samples $D_{corresponding}$ from the VGGSound dataset, and 7,500 wild data samples D_{random} from the Panda-70M dataset. Additionally, we randomly construct 70,000 audio-visual non-corresponding samples $D_{non-corresponding}$ based on VGGSound. We employ Imagebind (Girdhar et al., 2023) to extract and calculate the cosine similarity between the average representation of 8 random video frames and the audio representation. The similarity distribution curves of different sample sets are depicted in Figure 1. The audio-visual non-corresponding samples exhibit a normal distribution $N_{non-corresponding}(0.015, 0.081^2)$, while random wild samples follow the distribution $N_{random}(0.211, 0.116^2)$. In contrast, the audio-visual corresponding samples exhibit a left-skewed distribution with a higher concentration of similar instances. When the similarity of samples exceeds the threshold $\mu + 3\sigma$ (0.2564) of the audio-visual non-corresponding distribution $N_{non-corresponding}$, only 0.135% of the samples remain; thus, exceeding this threshold can be considered indicative of audio-visual correspondence. Notably, only 35% of the randomly selected wild data samples exhibit similarities exceeding the $\mu + 3\sigma$ (0.2564) threshold of the distribution $N_{non-corresponding}$.

Stage 3: Voice-Over Filtering. While the aforementioned filtering method effectively identifies non-corresponding samples based on audio-visual similarities, it fails to account for samples containing background music and voice-overs. These off-screen sounds, largely irrelevant to the visual content, can disrupt the intended audio-visual correspondence. To address this issue, we utilize the audio classification network PANNs (Kong et al., 2020) to label each audio clip, specifically targeting and filtering out these voice-overs. Following the classification scheme used in AudioSet, we annotate each audio clip with seven primary audio categories and their respective sub-categories.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

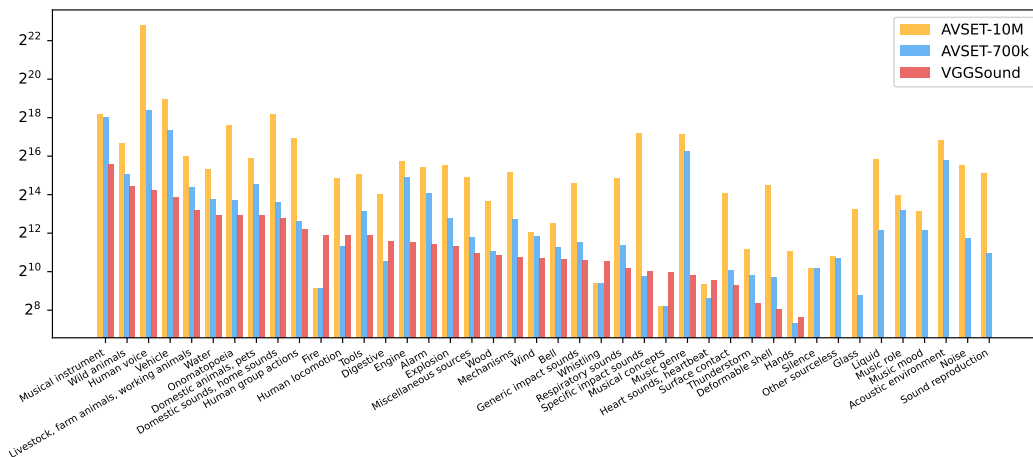


Figure 2: Comparison of the sample numbers for each audio category across AVSET-10M, AVSET-700K, and VGGSound datasets. Classification is carried out based on the secondary audio labels in AudioSet¹. We pseudo-label each sample from Panda-70M using PANNs (Kong et al., 2020), while labels on VGGSound are manually aligned with AudioSet.

Since speech and music are likely added during post-production, we specifically filter out samples that contain these elements along with other types of sounds. Other audio categories, such as the sounds of waterfalls and dog barking, typically originate from the original video. When these original video sounds co-occur with speech or music, it often indicates a high likelihood of off-screen voice interference. It is crucial to note that various instrumental sounds fall under the music category; thus, videos featuring instrumental performances are not excluded but are instead appropriately retained. Mirroring the approach in VGGSound (Chen et al., 2020), our filtering process aims to eliminate false positive samples—those with inappropriate sounds for each category. We refrain from using an audio classifier to select positive samples, as this may overlook some hard-to-classify yet criteria-meeting hard-positive audio samples.

Stage 4: Sample Recycling with Sound Separation. Speech is frequently encountered in wild video data, often leading to the inadvertent filtering out of a substantial amount of non-speech audio that includes voice-overs. This results in the loss of many potentially useful and valuable samples across various audio categories. Inspired by recent advancements in audio research (Jiang et al., 2023), we have implemented a sound separation model² (Solovyev et al., 2023) that is specifically designed to isolate sounds that are neither speech nor music from audio mixes contaminated with voice-over noise. The outputs from this sound separation process are subsequently returned to Stage 2 to verify the correspondence between the newly isolated audio and the video.

3.2 DATA ANALYSIS

We perform comprehensive statistical analyses on the AVSET-10M and AVSET-700K datasets to gain detailed insights. For further information about these datasets, please refer to Appendix D.

Diverse Categories, Abundant Samples. Figure 2 presents a comparative analysis of the number of audio categories in AVSET-10M, AVSET-700K, and VGGSound. To ensure consistency in audio category labels across different datasets, we employ the PANNs (Kong et al., 2020) audio classification network trained on AudioSet to label all samples in AVSET-10M. Subsequently, we manually align the labels in VGGSound with those in AudioSet and standardized the audio labels across all three datasets as secondary labels. It is evident that AVSET-10M and AVSET-700K encompass a broader range of audio types compared to VGGSound, including categories such as silence, liquid, and glass. Furthermore, AVSET-10M significantly outperforms AVSET-700K and VGGSound in most categories, offering a greater number of audio samples for each audio category.

¹<https://research.google.com/audioset/ontology/index.html>

²<https://github.com/ZFTurbo/MVSEP-CDX23-Cinematic-Sound-Demixing>

Table 2: Comparison of sample numbers after each stage. Due to partial video corruption, we could only download part of the original dataset. [†] The numbers here represent the video clips we collected. AVSET-10M (w/o. AVSET-700K) represents samples filtered from Panda-70M.

Stage	Goal	AVSET-700K		AVSET-10M (w/o. AVSET-700K)	
		#Num of Clips	Proportion	#Num of Clips	Proportion
<i>S1</i>	Candidate Videos [†]	1,445,360	100.0%	39,295,551	100.0%
<i>S2</i>	AV-C Filtering	898,366	62.2%	13,824,726	35.2%
<i>S3</i>	Voiceover Filtering	608,062	42.1%	7,124,923	18.1%
<i>S4</i>	Sample Recycling	727,530	50.3%	9,877,475	25.1%

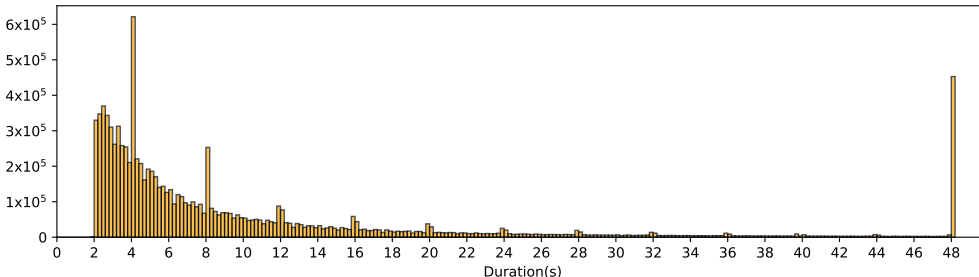


Figure 3: Histogram of Clip Length Distribution in AVSET-10M (w/o. AVSET-700K).

Duration Statistics. The samples filtered from Panda-70M include clips of varying lengths. As illustrated in Figure 3, we present the statistics for different clip lengths in AVSET-10M (excluding AVSET-700K). The total duration of AVSET-10M amounts to 30,418.6 hours, with an average clip length of 10.32 seconds. The longest clip spans 49 seconds, while the shortest measures 2 seconds. Notably, clips exceeding 10 seconds constitute 19,142.66 hours, representing 62.9% of total duration.

The Number of Video Samples after Each Filtering Stage. In Table 2, we detail the quantity of samples retained at each filtering stage for AVSET-700K and AVSET-10M (excluding AVSET-700K). Initially, in stage *S2* for AVSET-10M (excluding AVSET-700K), we filter out 64.8% of video samples due to lack of audio-visual correspondence. In the subsequent *S3* stage, 17.1% of the data containing voice-overs is removed. Further, in stage *S4*, an additional 8.0% of samples with voice-overs is refined through sound separation and subsequently recycled into the final audio-visual corresponding dataset. It is noteworthy that AudioSet undergoes a meticulous screening process by researchers, which results in a higher retention rate of data in the initial stage. AVSET-700K eliminates only 37.8% of data in its *S2* stage.

3.3 PRIVACY PROTECTION

All data in AVSET-10M was obtained through further screening of publicly available datasets (Gemmeke et al., 2017; Chen et al., 2024), with user permission obtained where necessary. In our work, we will only open-source the corresponding YouTube IDs and our annotations for these data samples, excluding any original data content. To further safeguard user privacy, we will implement a method that allows users to apply for the deletion of their corresponding samples. We will regularly synchronize user deletion requests with upstream datasets such as AudioSet and PANDA-70M to ensure compliance with privacy concerns.

3.4 DATASET VERIFICATION

We employ a distinct audio-visual representation learning model (Wang et al., 2024) different from the one used during the sample filtering phase to assess the reliability of our proposed sample filtering process. Specifically, we randomly sample data from four different audio-visual sources for validation: (1) audio-visual corresponding data from VGGSound, (2) audio-visual non-corresponding data created by randomly combining audio and video within VGGSound, (3) wild data randomly sampled from AudioSet, and (4) data from AVSET-700K obtained after the comprehensive filtering

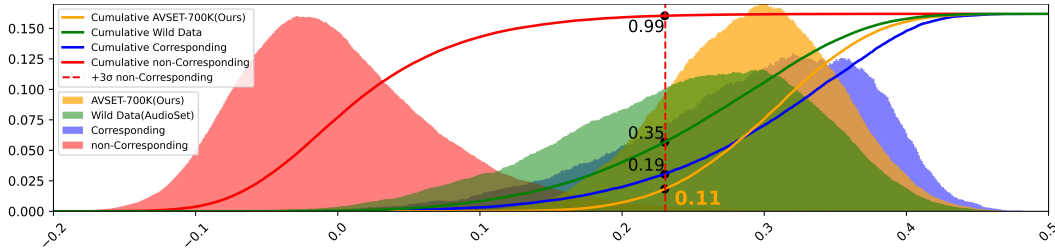


Figure 4: The distribution of audio-video cosine similarity of pre-trained model InternVL[†]_{1B}++(Ver.) Wang et al. (2024) was evaluated on different sample sets: (1) the audio-visual corresponding samples from VGGSound, (2) the randomly combined audio-visual non-corresponding samples from VGGSound, (3) the wild samples from AudioSet, and (4) the AVSET-700K sample set filtered with complete dataset processing. Notably, only 11% of the samples in AVSET-700K fall below the $\mu + 3\sigma$ threshold of non-corresponding distribution $N_{non-corresponding}$.

process. As depicted in Figure 4, we present the distributions of audio-visual similarity for these four sources. The mean and standard deviation of these similarities for each data source are detailed in Table 3.

AVSET-700K vs. AudioSet. It is evident that after data filtering, the audio-visual correspondence within the dataset is significantly enhanced compared to the wild data. The average cosine similarity of the AVSET-700K data increases from 0.258 to 0.303, while the standard deviation decreases from 0.086 to 0.058. Within the range $(\mu - 3\sigma, \mu + 3\sigma)$ of the normal distribution $N'_{non-corresponding}$ of non-corresponding data, the proportion of potential non-corresponding samples is reduced from 35% to 11%. This improvement demonstrates that our sample filtering method effectively enhances the audio-visual correspondence in the dataset.

Table 3: The mean and standard deviation (Std.) of audio-visual similarity among different sample sets.

Sample Sets	Mean	Std.
Non-Corresponding (Random)	0.015	0.072
Wild Data (AudioSet)	0.258	0.086
Corresponding (VGGSound)	0.302	0.083
AVSET-700K (ours)	0.303	0.058

AVSET-700K vs. VGGSound. As an audio-visual corresponding dataset, VGGSound contains a large number of samples with high audio-visual similarity. However, a substantial portion of the data exhibits low similarity, with 19% of VGGSound samples falling below the $\mu + 3\sigma = 0.231$ threshold of the distribution $N'_{non-corresponding}$. In contrast, only about 11% of the samples in AVSET-700K have an audio-visual similarity below 0.231, indicating that AVSET-700K contains more samples with high audio-visual correspondence. Additionally, AVSET-700K features a smaller standard deviation and fewer samples exhibiting extremely low similarity, demonstrating that our sample filtering process effectively enhances the robustness of audio-visual correspondence.

4 BENCHMARKS

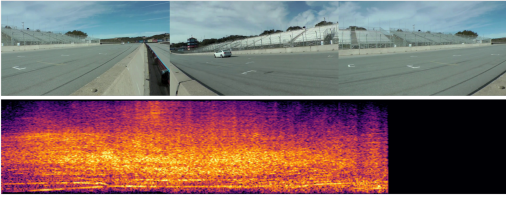
We benchmark two audio-visual tasks to explore the audio-visual correspondence: (1) Audio-Video Retrieval and (2) Vision-Queried Sound. In audio-video retrieval task, we experiment with AVSET-10M and focus on the data scale and the audio-visual temporally consistency. As for Vision-Queried Sound Separation, we mainly focus on the impact of each filtering stage, and work on the AVSET-700K which is of a similar scale to AudioSet. Specifically, we employ Imagebind (Girdhar et al., 2023) to extract the average features of 1 frame per second in the video as image features \mathbf{I} and InternVid (Wang et al., 2023) to extract the features of the entire video as video features \mathbf{V} . Please refer to Appendix C for additional details regarding the experiments.

4.1 AUDIO-VIDEO RETRIEVAL

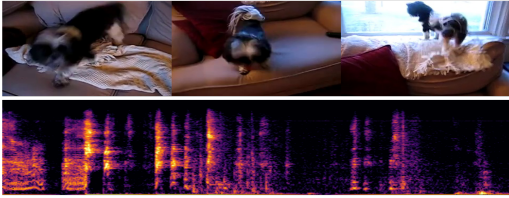
For the audio-video retrieval task, we validate on two audio-visual corresponding datasets, AVE (Tian et al., 2018) and VGGSound (Chen et al., 2020), and compare the Recall@1 (R@1) and Recall@5 (R@5) from vision to audio. For the image+video (I+V) modality, we apply feature weighting similar

Table 4: Comparison between the image-based method and the image+video based method on the task of visual to audio retrieval. The similarity on the diagonal should be the highest in each column. **The correct results** are highlighted in green, and **the incorrect results** are highlighted in red.


(a) Sample1 = $\{I_1, V_1, A_1\}$



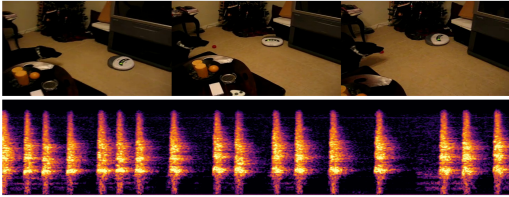
(b) Sample3 = $\{I_3, V_3, A_3\}$



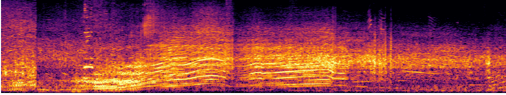
(c) Sample2 = $\{I_2, V_2, A_2\}$



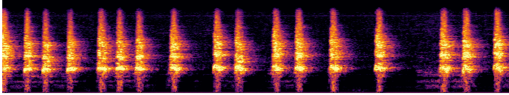
(d) Sample4 = $\{I_4, V_4, A_4\}$



(e) Similarity between Sample1 and Sample2.



(f) Similarity between Sample3 and Sample4.



I/V to A	I_1	I_2	I_1+V_1	I_2+V_2	I/V to A	I_3	I_4	I_3+V_3	I_4+V_4
A ₁	0.349	0.446	0.351	0.399	A ₃	0.373	0.416	0.388	0.304
A ₂	0.300	0.409	0.332	0.407	A ₄	0.402	0.457	0.357	0.359

Table 5: Comparison of vision to audio retrieval performance using different methods on ASE and VGGSound. **M** denotes the visual features used during retrieval.

#ID	M	Training Schedule	AVE		VGGSound	
			R@1	R@5	R@1	R@5
R1	I	AudioSet	18.00	40.11	11.74	28.52
R2	I	AVSET-700K	19.10	42.92	13.90	31.68
R3	I	AVSET-10M → AVSET-700K	19.11	43.05	13.91	31.94
R4	I+V	AVSET-700K	20.55	44.21	14.47	33.62
R5	I+V	AVSET-10M → AVSET-700K	20.78	44.47	14.93	34.03

to (Wang et al., 2024), with the mixed feature f_{I+V} calculated as $f_{I+V} = 0.9f_I + 0.1f_V$. In all the audio-video retrieval experiments conducted for this paper, we train a separate linear layer for each modality to align representations across different modalities, using a batch size of 1024.

AudioSet vs. AVSET-10M. AudioSet contains a significant number of audio-visual samples that do not correspond, adversely affecting audio-video alignment. By employing our filtered dataset, AVSET-700, we enhance cross-modal alignment capabilities, achieving a 3.16% improvement in VGGSound R@5 performance from R1 to R3 in Table 5. Furthermore, expanding the dataset to 10 million (R5) entries boosts the model performance on AVE R@5 by an additional 0.26%.

Based on Image vs. Based on Image+Video. Previous models, which rely solely on image features to retrieve audio clips that semantically match the image, lack the capability to evaluate audio-visual temporal consistency. As shown in Table 5, by leveraging both image and video features, the R@5 performance on VGGSound improved by 2.09% from R3 to R5, emphasizing the importance of audio-visual temporal consistency.

Qualitative Analysis. Table 4 presents several qualitative results of audio-video retrieval, underscoring the importance of temporal consistency for effective audio-video retrieval. For example, the image-based method could only deduce that engine roar should be present in the audio based on the image of a sports car, but it fails to determine when the sound should cease, leading to unsuccessful

Table 6: Comparison of sound separation performance among various methods on VGGSound. **M** stands for the query modality of sound separation.

#ID	M	Training Schedule	VGGSound	
			SDR \uparrow	SIR \uparrow
Baseline				
E1	I	VGGSound	5.606 \pm 0.102	8.074 \pm 0.161
E2	V	VGGSound	6.211\pm0.105	8.584\pm0.160
Zero-Shot				
E3	V	AudioSet	5.004 \pm 0.103	6.781 \pm 0.164
E4	V	AudioSet (w. AV-Correspondence Filtering)	5.646 \pm 0.101	7.682 \pm 0.162
E5	V	AVSET-700K	5.774\pm0.103	7.802\pm0.161
E6	V	AVSET-200K	5.152 \pm 0.103	6.928 \pm 0.168
Pretraining + Finetune				
E7	V	AudioSet (w. AV-Correspondence Filtering) \rightarrow VGGSound	6.548 \pm 0.103	9.251 \pm 0.158
E8	V	AVSET-700K \rightarrow VGGSound	6.666\pm0.102	9.377\pm0.158

audio-video pairing. In contrast, when both image and video features are considered, the similarity between mismatched sample pairs 1 and 2 is reduced from 0.446 to 0.399, thereby achieving correct audio-video pairing.

4.2 VISION-QUERIED SOUND SEPARATION

As shown in Table 6, we present the performance of vision-queried sound separation based on different modalities across various datasets. We utilize the framework of CLIPSep (Dong et al., 2022) to implement sound separation models across various modalities.

Image-Queried vs. Video-Queried. Compared to the sound separation model based on image queries (*E1*), the model utilizing video queries (*E2*) demonstrates superior performance, with the Signal-to-Distortion Ratio (SDR) improving by 0.605. This enhancement highlights the importance of audio-visual temporal consistency within the audio-visual research.

Corresponding vs. Non-Corresponding. Audio-visual correspondence is critical for effective sound separation. Models trained on the non-corresponding AudioSet (*E3*) encounter difficulties in achieving accurate separation and fail to capture proper audio-visual alignment. After implementing audio-visual correspondence filtering (*E4*), the dataset shows a marked improvement in audio-visual correspondence, as evidenced by a 0.642 increase in the Signal-to-Distortion Ratio (SDR). Despite this advancement, the presence of voice-over content continues to challenge the alignment between audio and visual modalities. Following a comprehensive filtering process, the model (*E5*) trained on AVSET-700K exhibits exceptional zero-shot sound separation capabilities, achieving an SDR of 5.774. This significant enhancement underscores the effectiveness of our proposed filtering process.

AVSET-200K vs. AVSET-700K. To further assess the impact of data scale on model performance, we randomly sampled 200K samples from AVSET-700K for experiments (*E6*). The performance dropped significantly, which demonstrates the importance of data scale. However, *E6* still outperformed *E3*, proving that audio-visual consistency is more critical than data scale.

5 CONCLUSION

Audio-visual correspondence datasets are pivotal for research in the audio-video domain. By applying a sample filtering process to AudioSet and Panda-70M, we have developed AVSET-10M—the first open, large-scale dataset with high audio-visual correspondence, comprising ten million audio-visual samples across 527 audio categories. Verification experiments demonstrate that AVSET-10M surpasses previous datasets in terms of audio-visual correspondence. Additionally, we benchmarked audio-video retrieval and vision-guided sound separation tasks, underscoring the critical role of audio-video temporal consistency in this field. Our AVSET-10M dataset opens up greater opportunities for advancement in audio-video research.

486 REPRODUCIBILITY STATEMENT

487

488 Our code has been open-sourced at <https://avset-10m.github.io/>, and the dataset will
 489 also be made publicly available upon acceptance. In Section 3.1, we provide a detailed explanation
 490 of the process for constructing the AVSET-10M dataset. Sections 4 and Appendix C outline the task
 491 definitions and specific implementation details, with the corresponding model training code also
 492 open-sourced.

493

494 REFERENCES

495

496 Pyscenedetect. In <https://github.com/Breakthrough/PySceneDetect>.

497

498 Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing
 499 Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text.
 500 *Advances in Neural Information Processing Systems*, 34:24206–24221, 2021.

501 Pranav Arora, Selen Pehlivan, and Jorma Laaksonen. Text-to-multimodal retrieval with bimodal
 502 input fusion in shared cross-modal transformer. In *Proceedings of the 2024 Joint International
 503 Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING
 504 2024)*, pp. 15823–15834, 2024.

505 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
 506 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
 507 *arXiv preprint arXiv:2308.12966*, 2023.

508

509 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
 510 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
 511 generation models as world simulators. 2024. URL [https://openai.com/research/
 512 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).

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

516

517 Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, Jinhui Tang, and Jing Liu.
 518 Valor: Vision-audio-language omni-perception pretraining model and dataset. *arXiv preprint
 519 arXiv:2304.08345*, 2023.

520 Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao,
 521 Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m:
 522 Captioning 70m videos with multiple cross-modality teachers. *arXiv preprint arXiv:2402.19479*,
 523 2024.

524

525 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
 526 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
 527 pp. 248–255. Ieee, 2009.

528 Hao-Wen Dong, Naoya Takahashi, Yuki Mitsufuji, Julian McAuley, and Taylor Berg-Kirkpatrick.
 529 Clipsep: Learning text-queried sound separation with noisy unlabeled videos. *arXiv preprint
 530 arXiv:2212.07065*, 2022.

531

532 Eduardo Fonseca, Jordi Pons Puig, Xavier Favory, Frederic Font Corbera, Dmitry Bogdanov, Andres
 533 Ferraro, Sergio Oramas, Alastair Porter, and Xavier Serra. Freesound datasets: a platform for the
 534 creation of open audio datasets. In *Hu X, Cunningham SJ, Turnbull D, Duan Z, editors. Proceedings
 535 of the 18th ISMIR Conference; 2017 oct 23-27; Suzhou, China.[Canada]: International Society
 536 for Music Information Retrieval; 2017. p. 486-93. International Society for Music Information
 537 Retrieval (ISMIR), 2017.*

538 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,
 539 Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*, 64(12):
 86–92, 2021.

- 540 Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing
541 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for
542 audio events. In *2017 IEEE international conference on acoustics, speech and signal processing*
543 *(ICASSP)*, pp. 776–780. IEEE, 2017.
- 544 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
545 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the*
546 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.
- 547 Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin
548 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced
549 diffusion models. In *International Conference on Machine Learning*, pp. 13916–13932. PMLR,
550 2023.
- 551 Vladimir Iashin and Esa Rahtu. A better use of audio-visual cues: Dense video captioning with
552 bi-modal transformer. *arXiv preprint arXiv:2005.08271*, 2020a.
- 553 Vladimir Iashin and Esa Rahtu. Multi-modal dense video captioning. In *Proceedings of the IEEE/CVF*
554 *conference on computer vision and pattern recognition workshops*, pp. 958–959, 2020b.
- 555 Sarah Ibrahim, Xiaohang Sun, Pichao Wang, Amanmeet Garg, Ashutosh Sanan, and Mohamed Omar.
556 Audio-enhanced text-to-video retrieval using text-conditioned feature alignment. In *Proceedings*
557 *of the IEEE/CVF International Conference on Computer Vision*, pp. 12054–12064, 2023.
- 558 Ziyue Jiang, Yi Ren, Zhenhui Ye, Jinglin Liu, Chen Zhang, Qian Yang, Shengpeng Ji, Rongjie
559 Huang, Chunfeng Wang, Xiang Yin, et al. Mega-tts: Zero-shot text-to-speech at scale with intrinsic
560 inductive bias. *arXiv preprint arXiv:2306.03509*, 2023.
- 561 Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung
562 Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on*
563 *Computer Vision and Pattern Recognition*, pp. 10124–10134, 2023.
- 564 Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig
565 Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model
566 for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.
- 567 Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. Panns:
568 Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Transactions*
569 *on Audio, Speech, and Language Processing*, 28:2880–2894, 2020.
- 570 Sangho Lee, Jiwan Chung, Youngjae Yu, Gunhee Kim, Thomas Breuel, Gal Chechik, and Yale
571 Song. Acav100m: Automatic curation of large-scale datasets for audio-visual video representation
572 learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
573 10274–10284, 2021.
- 574 Michael S Lew, Nicu Sebe, Chabane Djeraba, and Ramesh Jain. Content-based multimedia infor-
575 mation retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing,*
576 *Communications, and Applications (TOMM)*, 2(1):1–19, 2006.
- 577 Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer
578 questions in dynamic audio-visual scenarios. In *Proceedings of the IEEE/CVF Conference on*
579 *Computer Vision and Pattern Recognition*, pp. 19108–19118, 2022.
- 580 Wang Lin, Tao Jin, Wenwen Pan, Linjun Li, Xize Cheng, Ye Wang, and Zhou Zhao. Tadv: Towards
581 transferable audio-visual text generation. In *Proceedings of the 61st Annual Meeting of the*
582 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14983–14999, 2023.
- 583 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*
584 *neural information processing systems*, 36, 2024.
- 585 Simian Luo, Chuanhao Yan, Chenxu Hu, and Hang Zhao. Diff-foley: Synchronized video-to-audio
586 synthesis with latent diffusion models. *Advances in Neural Information Processing Systems*, 36,
587 2024.

- 594 Annamaria Mesaros, Aleksandr Diment, Benjamin Elizalde, Toni Heittola, Emmanuel Vincent,
595 Bhiksha Raj, and Tuomas Virtanen. Sound event detection in the dcase 2017 challenge. *IEEE/ACM*
596 *Transactions on Audio, Speech, and Language Processing*, 27(6):992–1006, 2019.
- 597
598 Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef
599 Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated
600 video clips. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.
601 2630–2640, 2019.
- 602 Tanzila Rahman, Bicheng Xu, and Leonid Sigal. Watch, listen and tell: Multi-modal weakly
603 supervised dense event captioning. In *Proceedings of the IEEE/CVF international conference on*
604 *computer vision*, pp. 8908–8917, 2019.
- 605 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
606 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
607 open large-scale dataset for training next generation image-text models. *Advances in Neural*
608 *Information Processing Systems*, 35:25278–25294, 2022.
- 609
610 Roman Solovyev, Alexander Stempkovskiy, and Tatiana Habruseva. Benchmarks and leaderboards
611 for sound demixing tasks, 2023.
- 612 Fabian-Robert Stöter, Antoine Liutkus, and Nobutaka Ito. The 2018 signal separation evaluation
613 campaign. In *Latent Variable Analysis and Signal Separation: 14th International Conference,*
614 *LVA/ICA 2018, Guildford, UK, July 2–5, 2018, Proceedings 14*, pp. 293–305. Springer, 2018.
- 615
616 Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland,
617 Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications*
618 *of the ACM*, 59(2):64–73, 2016.
- 619 Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization
620 in unconstrained videos. In *Proceedings of the European conference on computer vision (ECCV)*,
621 pp. 247–263, 2018.
- 622
623 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
624 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
625 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 626
627 Efthymios Tzinis, Scott Wisdom, Aren Jansen, Shawn Hershey, Tal Remez, Daniel PW Ellis, and
628 John R Hershey. Into the wild with audioscope: Unsupervised audio-visual separation of on-screen
629 sounds. *arXiv preprint arXiv:2011.01143*, 2020.
- 630
631 Efthymios Tzinis, Scott Wisdom, Tal Remez, and John R Hershey. Audioscopev2: Audio-visual
632 attention architectures for calibrated open-domain on-screen sound separation. In *European*
633 *Conference on Computer Vision*, pp. 368–385. Springer, 2022.
- 634
635 Yi Wang, Yanan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan
636 Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding
637 and generation. *arXiv preprint arXiv:2307.06942*, 2023.
- 638
639 Zehan Wang, Ziang Zhang, Xize Cheng, Rongjie Huang, Luping Liu, Zhenhui Ye, Haifeng Huang,
640 Yang Zhao, Tao Jin, Peng Gao, et al. Molecule-space: Free lunch in unified multimodal space via
641 knowledge fusion. *arXiv preprint arXiv:2405.04883*, 2024.
- 642
643 Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao,
644 Shibo Hao, Yemin Shi, Zhengzhong Liu, Eric P. Xing, and Zhiting Hu. Pandora: Towards general
645 world model with natural language actions and video states. 2024.
- 646
647 Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing:
Open-domain visual-audio generation with diffusion latent aligners, 2024.
- 648
649 Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and
650 Baining Guo. Advancing high-resolution video-language representation with large-scale video
651 transcriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
652 *Recognition*, pp. 5036–5045, 2022.

A LIMITATION

Since all upstream datasets of AVSET-10M rely on YouTube as the main data source, our dataset may be more closely aligned with the video styles prevalent on YouTube and may not fully represent video content from other platforms. However, to the best of our knowledge, our dataset is currently the largest audio-visual correspondence dataset available. In the future, we plan to verify the generalization ability of our AVSET-10M on data from other platforms. Additionally, we intend to collect data from a broader range of platforms to build a more diverse dataset.

B ETHICAL IMPACT

B.1 PRIVACY CONCERNS

AVSET-10M is built on existing open-source datasets and contains only video links, not the actual content. To address privacy concerns, we have implemented a deletion request mechanism that allows individuals to request the removal of links to privacy-sensitive content. Recognizing the limitations of users initiating such requests, we plan to periodically update our repository from upstream datasets (such as AudioSet and PANDA-70M) to proactively identify and remove any videos that may raise privacy concerns. This ensures continued adherence to privacy standards, as discussed in section 3.3.

B.2 POPULATION REPRESENTATIVENESS

Although privacy protection makes it challenging to determine the precise geographic location of videos, which complicates deep demographic analysis, we believe that the data samples offer a reasonable degree of population representativeness. Given that YouTube videos are uploaded by users all over the world, our dataset inherently captures a diverse range of demographics.

B.3 POTENTIAL APPLICATIONS

This paper primarily focuses on proposing a large-scale audio-visual correspondence dataset, aimed at expanding research possibilities in the audio-visual sector. This field includes technologies like video dubbing, which can lead to audio forgery. However, it's crucial to note that such dubbing does not result in severe identity forgery issues, unlike those caused by voice cloning technologies.

C IMPLEMENTATION DETAILS

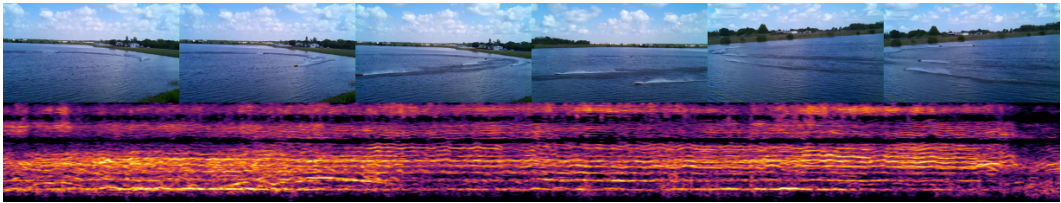
C.1 SOUND SEPARATION

Same as the experimental setting of (Dong et al., 2022), for all audio samples, we conduct experiments on samples of length 65535 (approximately 4 seconds) at a sampling rate of 16 kHz. For spectrum computation, we employ a short-time Fourier transform (STFT) with a filter length of 1024, a hop length of 256, and a window size of 1024. All images are resized to 224×224 pixels. All models are trained with a batch size of 128, using the Adam optimizer with parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$, for 200,000 steps. Additionally, we employ warm-up and gradient clipping strategies, following Dong et al. (2022). We compute the signal-to-distortion ratio (SDR) using museval (Stöter et al., 2018). All experiments are conducted on a single A800 GPU.

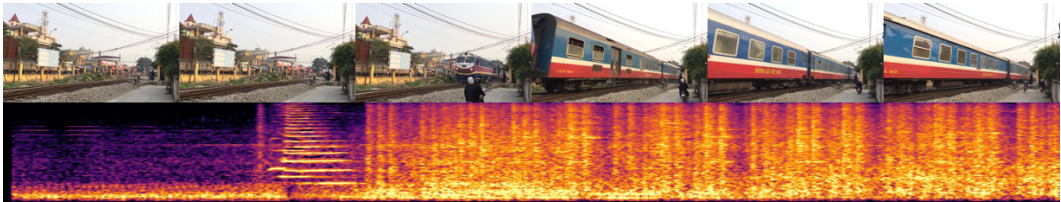
C.2 AUDIO-VIDEO RETRIEVAL

Same as the experimental setting of Wang et al. (2024), for all experiments, the softmax temperature is set to 0.01, and the temperature for the InfoNCE loss is set to 0.02. We utilize the Adam optimizer with a learning rate of 1×10^{-3} and a batch size of 2048, running the training process for 20 epochs.

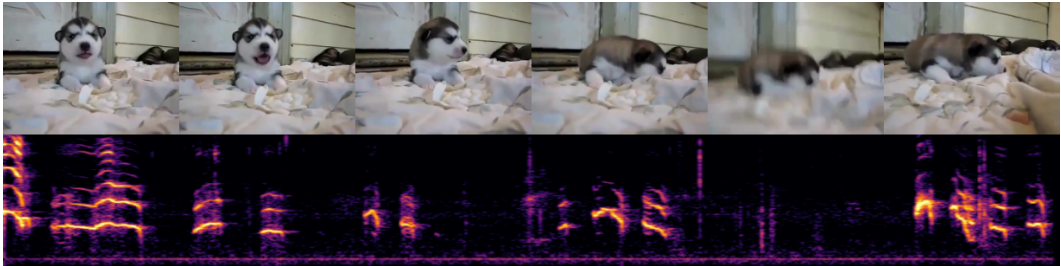
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755



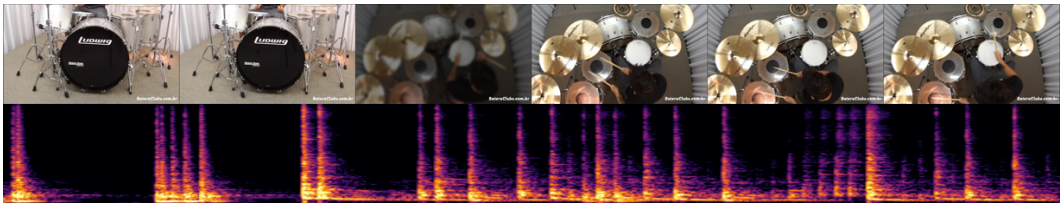
(a) Audio-Vision Cosine Similarity $\theta = 0.479$.



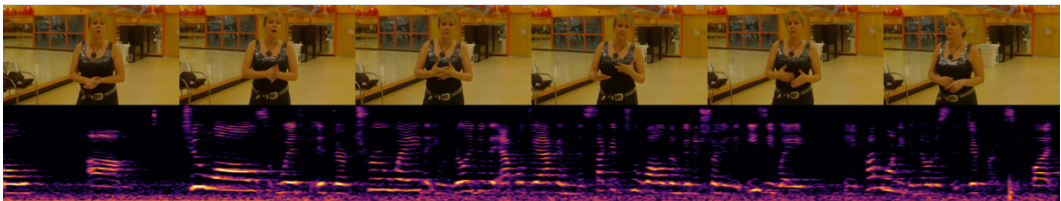
(b) Audio-Vision Cosine Similarity $\theta = 0.442$.



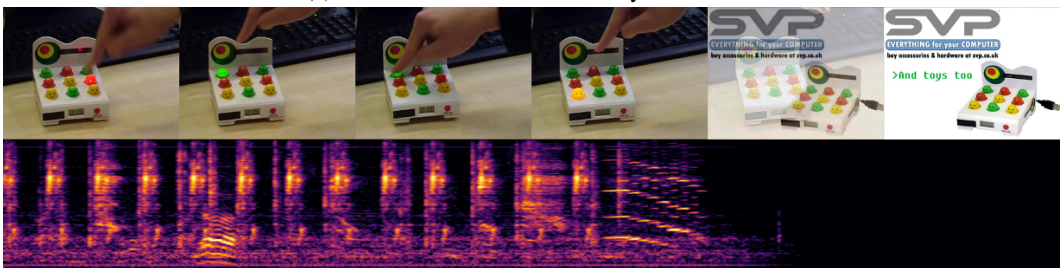
(c) Audio-Vision Cosine Similarity $\theta = 0.408$.



(d) Audio-Vision Cosine Similarity $\theta = 0.404$.



(e) Audio-Vision Cosine Similarity $\theta = 0.392$.



(f) Audio-Vision Cosine Similarity $\theta = 0.335$.

Figure 5: Audio-video consistency samples in AVSET.

756 D AVSET-10M

757 D.1 SAMPLES OF AVSET-10M

758 We present some audio-video consistency samples from the AVSET-10M in Figure 5. For additional
759 samples, please visit the demo page at <https://avset-10M.github.io>.

760 D.2 DATASET COMPOSITION

761 We release AVSET-10M as the following two subsets:

- 762 • **AVSET-700K**: This subset comprises 727,530 audio-visual corresponding samples filtered from
763 AudioSet. Each video segment in this subset is accompanied by a manually labeled audio category,
764 ensuring accurate categorization and relevance.
- 765 • **AVSET-10M (w/o. AVSET-700K)**: This subset comprises 10,234,280 audio-visual corresponding
766 samples, filtered from the Panda-70M dataset. Each video segment is semantically coherent,
767 focusing on a single event, and includes a text description originally from the Panda70M dataset.
768 Additionally, we provide pseudo-labels for the audio categories, derived with PANNs (Kong et al.,
769 2020), along with their corresponding confidence scores. Researchers can use these pseudo-labels
770 to freely partition the dataset.

771 We provide comprehensive meta-information for each video clip, including the YoutubeID of the
772 video, timestamps for each clip, audio-visual cosine similarity, a flag indicating whether sound
773 separation is required, and relevant text labels. For AVSET-10M (w/o. AVSET-700K), captions and
774 pseudo-labels are included, while AVSET-700K features manual audio labels.

775 D.3 LICENSE

776 AVSET-10M is released under the [CC BY 4.0] license. Before using this dataset, please ensure that
777 you have read and understood the terms of the license.

778 E DATASHEET OF AVSET-10M

779 We present a datasheet (Geburu et al., 2021) for documentation and responsible usage of LeanDojo
780 Benchmark.

781 E.1 MOTIVATION

- 782 1. **For what purpose was the dataset created?** We have developed the AVSET-10M dataset,
783 a tailored audio-video corresponding dataset, designed to advance audio-visual research by
784 facilitating the exploration of semantic and temporal alignment between audio and video
785 components.
- 786 2. **Who created the dataset and on behalf of which entity?** The AVSET-10M was developed
787 by researchers listed in the author list.

788 E.2 COMPOSITION

- 789 1. **What do the instance that comprise the dataset represent (e.g., documents, photos,
790 people, countries?)** Each instance consists of a pair of corresponding audio and video
791 samples, along with several associated labels.
- 792 2. **How many instances are there in total (of each type, if appropriate)?** The AVSET-10M
793 dataset contains 10,605,005 samples, of which the AVSET700K subset includes 727,530
794 samples.
- 795 3. **Does the dataset contain all possible instances or is it a sample of instances from a
796 larger set?** The dataset contains all possible instances.
- 797 4. **What data does each instance consist of?**

- 810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
5. **Is there a label or target associated with each instance?** We provide the cosine similarity between audio and visual content as well as the audio labels for each sample.
 6. **Is any information missing from individual instances?** For some instances filtered from Panda-70M, although the audio and video correspond, it is not able to identify the specific audio pseudo-labels. Note that this does not affect the audio-visual correspondence in our dataset.
 7. **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network links)?** N/A
 8. **Are there recommended data splits (e.g., training, development/validation, testing)?** In the AVSET-10M dataset, there are a large number of audio labels, allowing researchers to perform appropriate splits based on these labels. We do not have a recommended data splits.
 9. **Are there any errors, sources of noise, or redundancies in the dataset?** N/A
 10. **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?** We only provide the download links for the videos, the raw videos need to be downloaded from the YouTube platform.
 11. **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor–patient confidentiality, data that includes the content of individuals’ non-public communications)?** The AudioSet and Panda-70M used as the source contains facial videos that may pose a risk of infringement, we will delete the corresponding samples if necessary to avoid potential legal issues.
 12. **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?** Our data all come from the YouTube platform, which has a detailed data review process to ensure that it does not contain videos that are offensive, insulting, threatening, or might otherwise cause anxiety.
 13. **Does the dataset identify any subpopulations (e.g., by age, gender)?** N/A
 14. **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?** Individual identities may be identifiable through the video uploader.
 15. **Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)?** N/A

845 E.3 COLLECTION PROCESS

- 846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
1. **How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or language)?** The audio-video similarity is calculated using Imagebind (Girdhar et al., 2023), and the audio tags are obtained using PANNs (Kong et al., 2020).
 2. **What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?** All raw video data is sourced from established open-source datasets, and we employ an advanced filtering process to refine these data. The integrity and efficacy of the filtering process for the entire dataset have been thoroughly verified in Section 3.4.
 3. **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?** Based on the audio-video similarity.
 4. **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?** N/A.

- 864
865
866
867
868
5. **Were any ethical review processes conducted (e.g., by an institutional review board)?** Our data all come from the YouTube platform, which has a detailed data review process to ensure that it does not contain videos that are offensive, insulting, threatening, or might otherwise cause anxiety.

869 E.4 PREPROCESSING/CLEANING/LABELING

- 870
871
872
873
874
875
876
877
878
879
880
881
882
883
1. **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** We employ Imagebind (Girdhar et al., 2023) to determine the similarity between audio and video, PANNs (Kong et al., 2020) to classify audio into different categories, and a sound separation model (Solovyev et al., 2023) to extract non-speech tracks from the audio.
 2. **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?** We provide URLs for all raw videos, allowing researchers to download the videos directly from the YouTube platform.
 3. **Is the software that was used to preprocess/clean/label the data available?** ImageBind (<https://github.com/facebookresearch/ImageBind>). PANNs (https://github.com/qiuqiangkong/audioset_tagging_cnn). Sound Separation model (<https://github.com/ZFTurbo/MVSEP-CDX23-Cinematic-Sound-Demixing>).

884 E.5 USES

- 885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
1. **Has the dataset been used for any tasks already?** Yes, we have benchmarked the tasks of visual guided sound separation and audio-video retrieval using the AVSET-10M dataset.
 2. **Is there a repository that links to any or all papers or systems that use the dataset?** Yes. Please visit the web page of AVSET-10M (<https://avset-10M.github.io>).
 3. **What (other) tasks could the dataset be used for?** Our dataset is designed to facilitate research in video-to-audio generation, text-to-audio generation, and various other audio-video generation tasks. Additionally, it supports studies in audio-video classification, audio-video captioning, and other related audio-video understanding tasks.
 4. **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)?** To enlarge the sample size of non-speech categories, we utilize a sound separation model to process the data. This method may introduce a certain degree of audio distortion. Users can create a distortion-free sample set by using the identifiers provided in the dataset.
 5. **Are there tasks for which the dataset should not be used?** N/A.

903
904 E.6 DISTRIBUTION

- 905
906
907
908
909
910
911
912
913
914
915
1. **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** Yes, the dataset is open to the public.
 2. **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** The dataset will be distributed through platforms such as github and hugging face, and the code will be placed on github.
 3. **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** No.
 4. **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** No.

916 E.7 MAINTENANCE

- 917
1. **Who will be supporting/hosting/maintaining the dataset?** The first author of this paper.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

2. **Is there an erratum?** No. If errors are found in the future, we will release errata on the main web page for the dataset (<https://avset-10m.github.io/>).
3. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** Yes, the datasets will be updated whenever necessary to ensure accuracy, and announcements will be made accordingly. These updates will be posted on the main web page for the dataset (<https://avset-10m.github.io/>).
4. **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were the individuals in question told that their data would be retained for a fixed period of time and then deleted?)** The samples in the dataset are sourced from the YouTube platform. We have stated that if any specific fragments are found to infringe on individual rights, we will promptly remove them.
5. **Will older version of the dataset continue to be supported/hosted/maintained?** Yes, older versions of the dataset will continue to be maintained and hosted.
6. **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?** Our dataset will be published on the GitHub platform. If other researchers wish to further expand the dataset, they are welcome to contact us.