

ACAV-1M: DATA CURATION AND BENCHMARKING FOR AUDIO-VISUAL REPRESENTATION LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

The natural alignment of visual and audio information in videos provides a strong learning signal. However, commonly used large-scale video datasets contain audio-visual signals that are not aligned, e.g. background music. This limits the development of robust models that leverage the complementary nature of audio and video data. To address this limitation, we curate *ACAV-1M*, a new large-scale dataset that contains one million samples sourced from the *ACAV-100M* dataset. The *ACAV-1M* dataset is obtained through a pipeline that ensures the audio-visual correspondence and synchronization of samples in the dataset. Our pipeline transforms raw video and audio into text captions, followed by text summarization and an extensive filtering procedure. The filtering is done based on audio-caption alignment, audio-visual instance semantic alignment, and temporal synchronization. Furthermore, we propose an audio-visual learning benchmark that supports a diverse range of downstream tasks. Empirical evaluations demonstrate that models trained on *ACAV-1M* achieve superior performance compared to using existing datasets across all tasks. Our *ACAV-1M* dataset and code to reproduce all benchmark results will be made publicly available upon acceptance.

1 INTRODUCTION

Recent advancements in multimodal learning, exemplified by the Flan Collection (Longpre et al., 2023) and MMC4 (Zhu et al., 2023), have showcased significant strides for multimodal vision-language models. These developments underline the potential of integrated multimodal datasets for enhancing model performance. However, the field lacks a high-quality, large-scale collection specifically tailored for audio-visual learning, where audio and visual data complement each other to achieve a more holistic understanding of the environment.

The importance of a unified audio-visual dataset stems from the need for a systematic approach to evaluate the interaction between audio and visual inputs, which are often treated independently, as shown in Table 1. The integration of these modalities promises to improve the robustness and accuracy of learning models by leveraging their inherent complementary properties (Aytar et al., 2016; Owens et al., 2016; Arandjelovic & Zisserman, 2017; Korbar et al., 2018; Senocak et al., 2018). The absence of such datasets might hamper the development of audio-video models that can effectively exploit the synergies between sight and sound, thus limiting advancements in this area.

To address this, we introduce a new dataset, namely *ACAV-1M*, which consists of one million audio-visual samples. Our *ACAV-1M* follows a curation pipeline that consists of several steps. First, a multimodal Large Language Model (LLM) (Lin et al., 2023) generates multiple captions from audio and video inputs. Then, we use an LLM (OpenAI, 2023) to summarize the long captions into one sentence for the following quality measure steps. Lastly, we utilize ImageBind (Girdhar et al., 2023) to measure audio-language, audio-video instance, and audio-video temporal alignment for data curation. For audio-language alignment, we compute the normalized cosine similarity between audio instance features and caption features extracted from ImageBind (Girdhar et al., 2023). For audio-video instance alignment, we calculate the normalized cosine similarity between audio instance features and video features extracted from ImageBind. For audio-video temporal alignment, we compute the normalized cosine similarity between audio instance features and video features across all ten seconds extracted from ImageBind. These final alignment quality check steps ensure the coherence and synchronization between modalities.

Table 1: Details about dataset source, modality, number of samples, and benchmark tasks.

Dataset	Modality	# Data	Benchmark Tasks
ACAV100M (Lee et al., 2021)	Audio, Video	100M	Classification
AudioSet (Gemmeke et al., 2017)	Audio, Video	2.1M	Classification
Flickr-SoundNet (Aytar et al., 2016)	Audio, Video	2M	Classification, Localization
VGG-Sound (Chen et al., 2020b)	Audio, Video	200K	Classification, Localization
AudioCaps (Kim et al., 2019)	Audio, Video	48K	Retrieval
Kinetics-Sound (Arandjelović & Zisserman, 2017)	Audio, Video	19K	Classification
LLP (Tian et al., 2020)	Audio, Video	12K	Video Parsing
AVSD (AlAmri et al., 2019)	Audio, Video, Text	12K	Scene-Aware Dialog
MUSIC-AVQA (Li et al., 2022)	Audio, Video, Text	9K	Question-Answering
AVS-Bench (Zhou et al., 2022)	Audio, Video	7K	Segmentation
Clotho (Drossos et al., 2019)	Audio, Text	5K	Retrieval
AVE (Tian et al., 2018)	Audio, Video	4K	Localization
MUSIC (Zhao et al., 2018)	Audio, Video	448	Source Separation
ACAV-1M (ours)	Audio, Video, Text	1M	Cls. & SrcLoc. & Retrieval & SADialog. VideoPars. & QA & Seg. & SrcSep.

Our dataset and benchmark not only provide tools for measuring data quality on audio-visual instance alignment and temporal alignment, but also support an extensive range of downstream applications. These include audio-visual classification, sound source localization, retrieval, video parsing, scene-aware dialogue, audio-visual question-answering, segmentation, and sound source separation. We provide benchmark results for each of these applications with task-specific methods, and each of these applications is backed by benchmark baselines, task-specific methods, and both pre-trained and novel multimodal foundation models developed using *ACAV-1M*.

Empirical results from extensive experiments demonstrate that models trained on *ACAV-1M* surpass existing methods, highlighting the dataset’s effectiveness and scalability properties. This establishes *ACAV-1M* as a significant step towards the systematic integration of audio and visual data in machine learning research, providing a robust platform for exploring new frontiers in multimodal interaction and representation learning.

To summarize, we make the following four contributions:

- We curate the *ACAV-1M* dataset with one million audio-visual samples designed to address the gap in existing multimodal datasets for audio and visual data.
- Our data curation pipeline is a novel contribution that includes the transformation of raw video and audio into detailed, aligned captions using a multimodal large language model.
- We establish comprehensive benchmarks and task-specific methods that leverage our dataset to advance the state-of-the-art in audio-visual learning.
- Extensive experimental analyses demonstrate the effectiveness and scalability of models on *ACAV-1M* compared to existing audio-visual datasets.

2 RELATED WORK

Multimodal benchmarks. Dataset curation efforts, such as the Flan Collection (Longpre et al., 2023) and mmc4 (Zhu et al., 2023), have set precedents for multimodal learning. The Flan Collection has been instrumental in effective instruction tuning, while mmc4 addresses the challenges of few-shot, in-context, and interleaved learning across visual and language models. These benchmarks have laid the groundwork for *ACAV-1M* emphasizing the need for datasets that support intricate multimodal interactions.

Audio-visual learning. Audio-visual representations learning has been addressed in many previous works (Aytar et al., 2016; Owens et al., 2016; Arandjelovic & Zisserman, 2017; Korbar et al., 2018; Senocak et al., 2018; Zhao et al., 2018; 2019; Gan et al., 2020; Morgado et al., 2020; 2021a;b; Hershey & Casey, 2001; Ephrat et al., 2018; Hu et al., 2019). Exploiting the natural alignment across the audio and visual modalities is beneficial for many audio-visual tasks, such as audio-event localization (Tian et al., 2018; Lin et al., 2019; Wu et al., 2019; Lin & Wang, 2020), audio-visual localization (Morgado et al., 2018; Gao & Grauman, 2019; Chen et al., 2020a; Morgado et al., 2020), audio-visual navigation (Chen et al., 2020a; 2021a; 2022), and audio-visual parsing (Tian et al., 2020; Wu & Yang, 2021; Lin et al., 2021; Mo & Tian, 2022). Different to the aforementioned

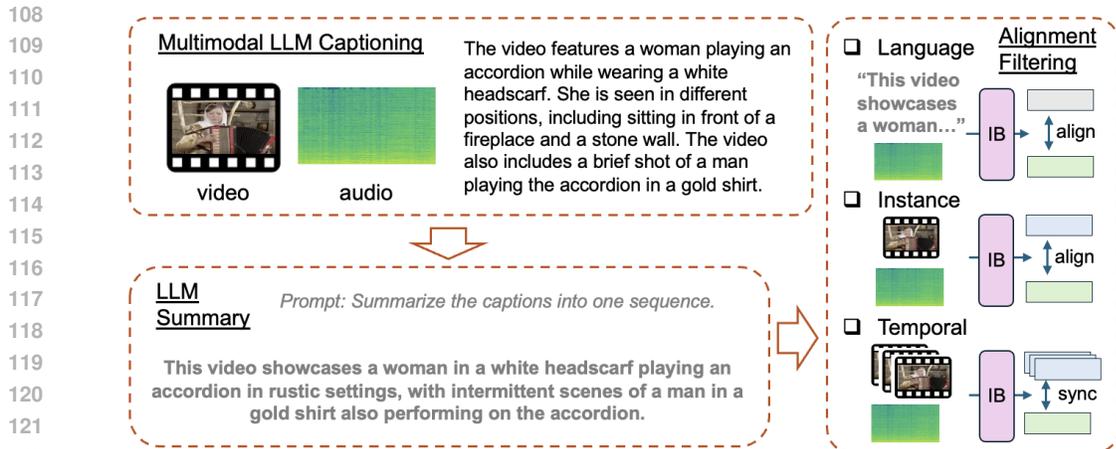


Figure 1: Illustration of the proposed *ACAV-1M* collection paradigm. First, we utilize a multimodal Large Language Model (LLM) (Lin et al., 2023; Mei et al., 2024) to generate multiple captions from video and audio inputs. Then, an LLM (OpenAI, 2023) is adopted to summarize the long captions into one sentence for the following alignment filtering steps. Finally, we use ImageBind (IB) (Girdhar et al., 2023) to ensure the coherence and synchronization between modalities by measuring audio-language, audio-video instance, and audio-video temporal alignment.

methods that focus on downstream applications, we propose a new dataset that gives boosts on downstream tasks when used for pre-training.

Audio-visual benchmarks. Existing audio-visual datasets, such as AudioSet (Gemmeke et al., 2017), VGGSound (Chen et al., 2020b), and ACAV100M (Lee et al., 2021), provide valuable resources for training and testing audio-visual models. These datasets have advanced audio-visual learning but are limited in size or contain noisy data or labels. Our *ACAV-1M* complements existing datasets by offering a structured and aligned dataset that facilitates cleaner multimodal integration.

3 *ACAV-1M* DATASET AND BENCHMARK

3.1 DATASETS CONSTRUCTION AND STATISTICS

ACAV-1M was meticulously constructed with a dataset curation process that ensures the close alignment between audio and visual elements which is crucial for complex multimodal learning tasks.

Data curation. The dataset curation follows a robust pipeline, starting with raw audio and video data, as illustrated in Figure 1. Each video clip is processed through our multimodal Large Language Model (LLM) to extract descriptive captions. Specifically, VideoLLaVA (Lin et al., 2023) and WavCaps (Mei et al., 2024) generate several sentence-level descriptions from video and audio inputs. These are condensed by a general LLM (OpenAI, 2023) into a single, comprehensive caption that captures the essence of the audio-visual content. This method ensures that our dataset supports semantic analysis and retrieval tasks effectively.

Data statistics. *ACAV-1M* is annotated with captions for both the audio and video components, facilitating cross-modal training and evaluation. Each audio segment is set to a 10-second duration to standardize the dataset and simplify the processing requirements. Audio data is categorized into various classes such as music, nature sounds, animal sounds, speech, machine noises, and others, providing a broad spectrum of audio types for comprehensive multimodal learning.

Alignment filtering. We quantify the alignment across modalities in our quality measure criteria. We employ ImageBind (Girdhar et al., 2023) to ensure several forms of alignment. 1) Language Alignment: The alignment between text captions and both audio and visual content is assessed with a normalized cosine similarity threshold of 0.5, ensuring that descriptions accurately reflect the content. 2) Instance Alignment: The synchronization between audio and visual streams is verified, with an emphasis on maintaining a normalized cosine similarity threshold of 0.5 to ensure alignment

162 across modalities. 3) Temporal Alignment: Audio and visual data are aligned within a temporal
 163 window of 1 second per segment, with an average alignment threshold of 0.5 across the dataset.
 164

165 **Algorithm 1** Pseudo Algorithm for Data Curation and Alignment Filtering
 166

167 1: **Input:** Raw video and audio data
 168 2: **Output:** Curated dataset with aligned audio-visual captions
 169
 170 3: **Data Curation Process:**
 171 4: **for** each video clip in dataset **do**
 172 5: Extract raw audio and video streams
 173 6: Use VideoLLaVA (Lin et al., 2023) to generate sentence-level descriptions from video
 174 7: Use WavCaps (Mei et al., 2024) to generate sentence-level descriptions from audio
 175 8: Condense descriptions using a general LLM (OpenAI, 2023) into a single comprehensive
 176 caption
 177 9: Attach the comprehensive caption to the corresponding video clip
 178 10: **end for**
 179
 180 11: **Alignment Filtering Process:**
 181 12: **for** each item in curated dataset **do**
 182 13: **Language Alignment:**
 183 14: Calculate normalized cosine similarity between text captions and audio-visual content
 184 15: **if** similarity < 0.5 **then**
 185 16: Flag for review or reprocessing
 186 17: **end if**
 187 18: **Instance Alignment:**
 188 19: Assess synchronization between audio and visual streams using ImageBind (Girdhar et al.,
 189 2023)
 190 20: **if** similarity < 0.5 **then**
 191 21: Adjust synchronization parameters and re-align
 192 22: **end if**
 193 23: **Temporal Alignment:**
 194 24: Check for alignment within a temporal window of 1 second per segment
 195 25: **if** average alignment threshold < 0.5 **then**
 196 26: Refine temporal synchronization parameters
 197 27: **end if**
 198 28: **end for**
 199
 200 29: **Return** finalized dataset with validated and aligned captions

201 **Algorithm for ACAV-1M Data Curation and Alignment.** Algorithm 1 is a pseudo-algorithm that
 202 encapsulates the data curation and alignment filtering processes described for the ACAV-1M dataset.
 203 This algorithm is structured to provide a clear, step-by-step procedure that reflects the robust method-
 204 ologies used in preparing the dataset. This algorithm also provides a structured approach to process-
 205 ing and aligning the data within the ACAV-1M ensuring that each component (video, audio, and
 206 textual caption) is effectively synchronized and semantically coherent. The algorithm is designed to
 be part of a larger document or paper, offering clarity on the methods and steps taken to curate and
 align data within the dataset.

207 3.2 AUDIO-VISUAL BENCHMARK

208 **Audio-visual tasks.** ACAV-1M supports an extensive range of audio-visual downstream tasks, each
 209 designed to leverage the rich, multimodal nature of the ACAV-1M dataset.
 210

- 211 1. **Audio-Visual Classification.** The goal is to classify the scenes or objects depicted in the audio-
 212 visual clips with accuracy as evaluation. For linear probing and fine-tuning on audio-visual clas-
 213 sification, we used VGGSound-Music with 49 classes and VGGSound-All with 221 categories.
- 214 2. **Audio-Visual Source Localization.** This task measures the model’s ability to localize sound
 215 sources within a visual frame, assessed by the mean Intersection over Union (mIoU). We use
 Flickr-SoundNet (Senocak et al., 2018) with 4,500 pairs for training and testing the model on

- 216 250 audio-visual pairs of sounding objects and extended 250 non-sounding objects introduced in
 217 SLAVC (Mo & Morgado, 2022b).
- 218 3. **Audio-Visual Retrieval.** The focus is the recall of relevant audio-visual content based on query
 219 descriptions. We use MSR-VTT (Xu et al., 2016) that includes 10K YouTube videos with 200K
 220 description sentences, where 9K is split for training and 1K for testing.
 - 221 4. **Audio-Visual Scene-Aware Dialog.** This task focuses on generating dialogues that are contextually
 222 relevant to given audio-visual scenes, evaluated via BLEU and METEOR scores. We use the
 223 AVSD track of the 10-th Dialog System Technology Challenges (DSTC10) (Shah et al., 2022)
 224 dataset.
 - 225 5. **Audio-Visual Video Parsing.** This involves parsing complex video scenes into simpler seg-
 226 ments, evaluated using an F-score at a mIoU threshold of 0.5. The LLP dataset (Tian et al., 2020)
 227 contains 11,849 YouTube video clips of 10-seconds long from 25 different event categories, such
 228 as car, music, cheering, speech, etc. We follow the official splits (Tian et al., 2020) of validation
 229 and test sets to train and test.
 - 230 6. **Audio-Visual Question-Answering.** This task tests accuracy in answering questions based on
 231 the content depicted in the audio-visual clips. we use the MUSIC-AVQA (Li et al., 2022) dataset
 232 that consists of 45,867 question-answer pairs and 9,288 videos.
 - 233 7. **Audio-Visual Segmentation.** This task focuses on the segmentation masks of visual elements,
 234 with performance measured by the F1 Score. AVSBench (Zhou et al., 2022) includes 4,932
 235 videos (in total 10,852 frames) from 23 categories, including instruments, humans, animals, etc.
 236 We use the official split of 3,452/740/740 videos for train/val/test.
 - 237 8. **Audio-Visual Source Separation.** The objective is to measure the ability to isolate individual
 238 audio sources from a mixed audio track, evaluated using metrics such as Signal to Distortion
 239 Ratio (SDR), Signal to Interference Ratio (SIR), and Signal to Artifacts Ratio (SAR). We use
 240 VGGSound-Music (Mo & Morgado, 2023) with 40,908 video clips from 49 music categories for
 241 training and 1201 clips for testing. VGGSound-Instruments (Hu et al., 2022) includes 32k video
 242 clips of 10-second length from 36 musical instrument classes, a subset of VGG-Sound (Chen
 243 et al., 2020b), and each video only has one single instrument class annotation. MUSIC (Zhao
 244 et al., 2018) consists of 448 untrimmed YouTube music videos of solos and duets from 11 instru-
 245 ment categories, where we use 358 solo videos for training and 90 solo videos for evaluation.

244 Here, we explain the baselines, task-specific methods, pre-trained models, and multimodal founda-
 245 tion used in our audio-visual benchmark.

246 **Task-specific methods.** *ACAV-1M* is utilized to establish a variety of task-specific methods tailored
 247 to each downstream task. For Audio-Visual Classification, methods are optimized for maximum ac-
 248 curacy. In Audio-Visual Source Localization, algorithms focus on improving the mean Intersection
 249 over Union (mIoU). For Audio-Visual Retrieval, the emphasis is on enhancing recall rates. Sim-
 250 ilarly, task-specific approaches are devised for Audio-Visual Video Parsing, Scene-Aware Dialog,
 251 Question-Answering, Segmentation, and Source Separation, each aiming to excel in metrics such as
 252 F-score, BLEU, METEOR, and Signal Decomposition Ratings (SDR, SIR, SAR).

253 **Pre-trained models.** We evaluate several models pre-trained on *ACAV-1M* including audio-
 254 MAE (Huang et al., 2022b), CAV-MAE (Gong et al., 2023), MAViL (Huang et al., 2022a), and
 255 AVMAE (Georgescu et al., 2023). These models leverage masked autoencoding techniques tailored
 256 for either audio alone or audio-visual data on audio-visual classification to have a comprehensive
 257 understanding on these models.

258 **Our method.** We use audio-visual masked autoencoders (He et al., 2021; Huang et al., 2022b)
 259 with masked modeling objectives. Specifically, we apply a modality-specific encoder with self-
 260 attention transformers to encode unmasked patches and use a decoder to predict the masked patches
 261 of the input modality from unmasked encoded and masked tokens. The overall model is simply
 262 optimized to reconstruct the original input modality of masked tokens using a ℓ_2 norm objective
 263 across predicted audio/visual tokens $\hat{\mathbf{x}}_m^a, \hat{\mathbf{x}}_m^v$ and ground-truth tokens $\mathbf{x}_m^a, \mathbf{x}_m^v$ defined as:

$$264 \mathcal{L} = \frac{1}{M^a} \sum_{m=1}^{M^a} \|\mathbf{x}_m^a - \hat{\mathbf{x}}_m^a\|_2^2 + \frac{1}{M^v} \sum_{m=1}^{M^v} \|\mathbf{x}_m^v - \hat{\mathbf{x}}_m^v\|_2^2, \quad (1)$$

265 where M^a, M^v denote sets of random masks applied on the input patch embeddings for audio and
 266 visual tokens, separately.

Table 2: **Audio-visual classification** on the VGGSound-Music, VGGSound-All, and AudioSet.

Method	VGGSound-Music		VGGSound-All		AudioSet	
	Linear (%)	Finetune (%)	Linear (%)	Finetune (%)	Linear (%)	Finetune (%)
MAE (He et al., 2021)	25.32	52.39	15.61	45.73	11.52	24.23
AudioMAE (Huang et al., 2022b)	41.65	55.61	42.35	57.76	30.23	44.92
CAV-MAE (Gong et al., 2023)	60.53	67.26	55.27	65.53	40.56	51.29
MAViL (Huang et al., 2022a)	61.95	69.53	57.36	67.17	43.62	53.38
AV-MAE (Georgescu et al., 2023)	60.82	67.61	56.15	65.08	41.67	51.32
<i>ACAV-1M</i> (ours)	64.87	71.25	61.35	69.29	47.83	56.05

Table 3: **Audio-visual source localization.** Quantitative results on Flickr-SoundNet.

Method	Precision	AP	F1
Attention 10k (Senocak et al., 2018)	49.38	51.23	55.39
OTS (Arandjelovic & Zisserman, 2018)	51.23	53.28	58.12
DMC (Hu et al., 2019)	50.52	52.93	57.56
CoarsetoFine (Qian et al., 2020)	51.76	54.85	58.63
DSOL (Hu et al., 2020)	55.29	57.92	62.05
LVS (Chen et al., 2021b)	52.38	55.31	59.35
EZVSL (Mo & Morgado, 2022a)	54.71	57.51	61.38
Mix-and-Localize (Hu et al., 2022)	55.83	58.21	62.52
SLAVC (Mo & Morgado, 2022b)	55.65	58.12	62.39
<i>ACAV-1M</i> (ours)	58.67	60.75	65.02

Table 4: **Audio-video retrieval.** Quantitative results on the MSR-VTT dataset.

Method	R@1	R@5	R@10
AVLnet (Rouditchenko et al., 2020)	19.62	50.32	60.51
TVLT (Tang et al., 2022)	23.83	52.56	63.92
<i>ACAV-1M</i> (ours)	26.57	58.78	70.26

Table 5: **Audio-visual scene-aware dialog.** Quantitative results on the DSTC10 dataset.

Method	BLEU	METEOR
MMA (Hori et al., 2018)	24.91	19.36
BMT (Iashin & Rahtu, 2020)	36.23	22.83
JST (Shah et al., 2022)	38.52	24.71
<i>ACAV-1M</i> (ours)	43.27	28.65

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. We use audio-visual pairs from our *ACAV-1M* dataset for pre-training. We finetune the model on datasets specific to the downstream tasks, as described in Section 3.2.

Evaluation metrics. Following the prior work (Hu et al., 2022; Mo & Morgado, 2022a;b), we use the Precision and F1 scores defined in (Mo & Morgado, 2022b) for visual source localization. For source separation, following (Zhao et al., 2018), we use Signal-to-Distortion Ratio (SDR) and Signal-to-Artifact Ratio (SAR). For audio-visual segmentation, we apply mIoU and F1 scores as evaluation metrics, following the previous work (Zhou et al., 2022). Linear probing and fine-tuning classification evaluations are based on top-1 accuracy, which measures the class difference from the ground-truth labels. For video parsing, we use F-scores to evaluate segment-level predictions for audio-visual events and Type@AV & Event@AV for the overall evaluation performance.

Implementation. The input images are resized to 224×224 . The audio is represented by log spectrograms extracted from 10s of audio at a sample rate of 8000Hz. We follow the prior work (Mo & Morgado, 2022a) and apply STFT to generate an input tensor of size 128×128 (128 frequency bands over 128 timesteps) using 50ms windows with a hop size of 25ms. For the audio and visual encoder, we use single-modality MAEs (He et al., 2021; Huang et al., 2022b). The models were trained on four A100 GPUs for 100 epochs using the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $1e - 4$ and a batch size of 128.

4.2 BENCHMARK EXPERIMENTAL RESULTS

Audio-Visual Classification. To validate the effectiveness of *ACAV-1M* on audio-visual classification, we compare to the following prior baselines: 1) MAE (He et al., 2021): a masked autoencoder with only images as input; 2) AudioMAE (Huang et al., 2022b): a masked autoencoder with only audio as input; 2) Audio-Visual MAEs (Gong et al., 2023; Huang et al., 2022a; Georgescu et al., 2023): masked autoencoders with both audio and images as input. Table 2 reports the quantitative comparison results. On VGGSound-Music, we achieved top results with 64.87% in linear probing and 71.25% in fine-tuning, indicating robustness in music-specific scenes. For VGGSound-All, we also recorded 61.35% in linear probing and 69.29% in fine-tuning, showcasing versatility across diverse audio-visual contexts. Regarding AudioSet, our model performed well with 47.83% in linear probing and 56.05% in fine-tuning, reflecting strong generalization capabilities.

Audio-Visual Source Localization. To validate the effectiveness of the proposed *ACAV-1M* dataset for sound source localization, we compare to the following prior work: 1) Attention 10k (Senocak et al., 2018) (CVPR 2018): the first baseline on sound source localization using a two-stream and attention-based neural network; 2) OTS (Arandjelovic & Zisserman, 2018) (ECCV 2018): a correspondence-based baseline for localization; 3) DMC (Hu et al., 2019) (CVPR 2019): a deep multi-modal clustering approach based on audio-visual co-occurrences; 4) CoarsetoFine (Qian et al., 2020) (ECCV 2020): a two-stage approach using coarse-to-fine embedding alignment; 5) DSOL (Hu et al., 2020) (NeurIPS 2020): a class-based method with two-stage training; 6) LVS (Chen et al., 2021b) (CVPR 2021): a contrastive learning framework with hard negative mining to learn audio-visual correspondence maps; 7) EZ-VSL (Mo & Morgado, 2022a) (ECCV 2022): a recent weakly supervised localization framework based on multiple-instance contrastive learning; 8) Mix-and-Localize (Hu et al., 2022) (CVPR 2022): a recent method based on a contrastive random walk on a graph of images and separated sound sources. 9) SLAVC (Mo & Morgado, 2022b) (NeurIPS 2022): a strong baseline with momentum encoders and extreme visual dropout to identify negatives and solve significant overfitting. The results are reported in Table 3. As can be seen, our *ACAV-1M* scored 58.67% Precision, which is the highest among the compared methods, indicating a high accuracy in predicting the correct localization of sound sources. We also achieved 60.75% Average Precision (AP), highlighting the method’s consistent performance across different thresholds, outperforming other methods in handling diverse scenarios. Our model achieves a 65.02% F1 score, which reflects the balance between precision and recall, demonstrating the robustness of our approach to effectively localize sound sources.

Audio-Visual Retrieval. For audio-visual retrieval, we evaluated the performance of our *ACAV-1M* model against established methodologies. This evaluation was performed using the MSR-VTT dataset, a comprehensive and challenging benchmark for video understanding and retrieval tasks, where we compare to the following baselines: 1) AVLnet (Rouditchenko et al., 2020): A self-supervised learning approach that develops a joint audio-visual-textual embedding space, leveraging the natural synchrony in videos to align raw video, audio, and text signals without requiring manual annotations. 2) TVLT (Tang et al., 2022): A very recent approach that introduces a visual-audio pre-training framework and incorporates masked audio/video autoencoding coupled with contrastive modeling, which aims to fine-tune the alignment between video and audio modalities to improve retrieval accuracy. The experimental results are shown in Table 4. In particular, we achieved 26.57% R@1, significantly higher than AVLnet (19.62% R@1) and TVLT (23.83% R@1), indicating a more precise retrieval at the topmost rank. Meanwhile, we scored 58.78% R@5, surpassing both AVLnet (50.32% R@5) and TVLT (52.56% R@5). Our *ACAV-1M* model demonstrates superior performance across all recall metrics.

Audio-Visual Scene-Aware Dialog. In the task of audio-visual scene-aware dialog, model are evaluated to demonstrate their capability to generate contextually appropriate dialog based on both visual and auditory inputs. We compared the performance against several prominent methods in the field: 1) MMA (Hori et al., 2018): an end-to-end conversation model that generates dialog responses based on multimodal attention-based video features, which integrates audio and visual cues to form a comprehensive understanding of the video content. 2) BMT (Iashin & Rahtu, 2020): a bi-modal Transformer that adapts the traditional Transformer architecture for bi-modal inputs, processing both audio and visual modalities to enhance performance on tasks like dense video captioning. 3) JSTL (Shah et al., 2022): a recent AV-transformer that employs attentional multimodal fusion and combines joint student-teacher learning and model combination techniques to refine dialog generation based on audio-visual data. Table 5 reports the results on the DSTC10 dataset. We observe a BLEU score of 43.27, surpassing all other compared models and reflecting its superior ability to generate grammatically and semantically correct sentences. With a METEOR score of 28.65, our model also leads in this metric.

Audio-Visual Video Parsing. In audio-visual video parsing, we conducted a comparative analysis using the LLP dataset. The comparison the following approaches: 1) AVE (Tian et al., 2018): An audio-guided co-attention network which includes additional branches for audio-visual parsing. This model leverages audio cues to enhance the segmentation and identification of visual elements in video. 2) AVSDN (Lin et al., 2019): A dual sequence-to-sequence model that merges global audio-visual features into localized contexts. This model aims to improve the parsing accuracy by enhancing the interaction between audio and visual modalities. 3) HAN (Tian et al., 2020): A hybrid attention network that utilizes multimodal multiple instance learning pooling. This network focuses on capturing the intricate relationships between audio and visual cues within video content to refine parsing accuracy. 4) MGN (Mo & Tian, 2022): A Multi-modal Grouping Network that aggregates

Table 6: **Audio-visual video parsing.** Quantitative results on the LLP dataset.

Method	Audio-Visual	Type@AV	Event@AV
AVE (Tian et al., 2018)	35.43	39.92	41.63
AVSDN (Lin et al., 2019)	37.12	45.73	50.82
HAN (Tian et al., 2020)	48.92	54.03	55.42
MGN (Mo & Tian, 2022)	50.63	55.62	57.25
<i>ACAV-1M</i> (ours)	55.35	58.96	58.67

Table 7: **Audio-visual question answering.** Quantitative results on MUSIC-AVQA.

Method	A	V	AV
AVSD (Schwartz et al., 2019)	68.52	70.83	65.49
Pano-AVQA (Yun et al., 2021)	70.73	72.56	66.64
AVQA (Li et al., 2022)	74.06	74.00	69.54
<i>ACAV-1M</i> (ours)	76.87	76.65	73.25

Table 8: **Audio-visual segmentation.** Quantitative results on the AVSBench dataset.

Method	mIoU	F1
Attention 10k (Senocak et al., 2018)	20.76	31.25
OTS (Arandjelovic & Zisserman, 2018)	24.55	36.85
DMC (Hu et al., 2019)	23.51	35.27
CoarsetoFine (Qian et al., 2020)	26.53	38.62
DSOL (Hu et al., 2020)	29.85	42.23
LVS (Chen et al., 2021b)	27.32	40.18
EZVSL (Mo & Morgado, 2022a)	30.52	43.26
Mix-and-Localize (Hu et al., 2022)	31.69	45.35
SLAVC (Mo & Morgado, 2022b)	31.36	45.02
<i>ACAV-1M</i> (ours)	36.39	49.85

event-aware unimodal features through semantically-aware grouping. It employs learnable categorical embedding tokens. Table 6 shows the experimental results. We achieved an accuracy of 55.35 and a Type@AV score of 58.96, the highest among all compared models, showcasing its exceptional ability to classify and understand different types of content accurately. Furthermore, we scored an Event@AV score of 58.67, illustrating strong performance in identifying and segmenting specific events within videos.

Audio-Visual Question-Answering. For audio-visual question answering (AVQA), our model was assessed on the MUSIC-AVQA dataset, testing its capability to integrate and interpret audio, visual, and combined audio-visual information to answer related questions accurately. This performance was benchmarked against the following models: 1) AVSD (Schwartz et al., 2019): a straightforward approach for audio-visual scene-aware dialog, trained end-to-end to tackle AVQA by directly associating audio-visual scenes with dialog responses. 2) Pano-AVQA (Yun et al., 2021): a multimodal transformer encoding with a unique approach to attention mechanisms that incorporate both audio and visual inputs simultaneously. 3) AVQA (Li et al., 2022): a very recent baseline that integrates comprehensive multimodal information by associating spatial grounding, temporal grounding, and advanced multimodal fusion techniques. Table 7 illustrates the experimental results on the MUSIC-AVQA dataset. For instance, we achieved a 76.87%@Audio score, indicating a high proficiency in extracting and utilizing audio information to answer questions, outperforming all other models. With a score of 73.25%@Audio-Visual, our model demonstrates a superior ability to use audio and visual data for answering questions, surpassing other methodologies in effectively utilizing integrated multimodal cues.

Audio-Visual Segmentation. In audio-visual segmentation, we did a comparative analysis using the AVSBench (Zhou et al., 2022) dataset, which is designed to evaluate segmentation capabilities across models that integrate audio and visual data. This task extends beyond localization to include the generation of accurate segmentation masks for audio-visual sources. We use the same baselines (Senocak et al., 2018; Arandjelovic & Zisserman, 2018; Hu et al., 2019; Qian et al., 2020; Hu et al., 2020; Chen et al., 2021b; Mo & Morgado, 2022a; Hu et al., 2022; Mo & Morgado, 2022b) as those for audio-visual source localization, adapted to generate detailed segmentation masks rather than just coarse localization maps. The results are reported in Table 8. Our model achieved a mIoU score of 36.39, indicating superior accuracy in segmenting relevant audio-visual content precisely. We also achieved an F1 score of 49.85, the highest among all compared methods. These results highlight the *ACAV-1M* model’s robust capability to accurately segment complex audio-visual scenes, establishing it as a leading method for audio-visual segmentation.

Audio-Visual Source Separation. To demonstrate the effectiveness of the proposed *ACAV-1M* on source separation, we compare to the following methods: 1) NMF (Virtanen, 2007): a traditional signal processing approach based on non-negative matrix factorization to generate the spectrogram of each sound source; 2) RPCA (Huang et al., 2012): a parameter-free baseline based on robust principal component analysis; 3) Sound-of-Pixels (Zhao et al., 2018): a deep learning approach that recovers separated audio conditioned on pixel-level visual features; 4) MP-Net (Xu et al., 2019): an improved audio-visual method based on recursive separation from the mixture; 5) CCoL (Tian et al., 2021) (CVPR 2021): a cyclic co-learning framework based on sounding object visual grounding to separate individual sound sources. 6) OneAVM (Mo & Morgado, 2023) (ICML 2023): a unified audio-visual framework for localization, separation, and recognition. We report the compar-

Table 9: **Sound source separation.** Quantitative results on the MUSIC and VGGSound datasets.

Method	MUSIC		VGGs-Instruments		VGGs-Music	
	SDR	SAR	SDR	SAR	SDR	SAR
NMF (Virtanen, 2007)	-0.62	2.41	-3.85	-0.76	-7.12	-9.01
RPCA (Huang et al., 2012)	0.86	3.81	-2.39	1.58	-5.53	-7.82
Sound-of-Pixels (Zhao et al., 2018)	4.55	10.24	2.52	4.67	0.95	1.03
MP-Net (Xu et al., 2019)	4.82	10.56	2.63	4.85	1.37	1.39
CCoL (Tian et al., 2021)	6.35	9.75	3.28	5.01	2.07	2.18
OneAVM (Mo & Morgado, 2023)	7.38	7.48	5.36	5.52	2.51	2.61
<i>ACAV-1M</i> (ours)	10.75	11.23	8.23	8.38	5.06	5.32

Table 10: **Ablation results on the benefit of our data curation pipeline** across different audio-visual benchmarks. Note that we use 100K VGGSound samples for pre-training.

Data Curation	Cls. Acc (%)	SrcLoc. Prec	Retrieval Acc (%)	SADialog. BLEU	VidPars. F-score (%)	QA Acc (%)	Seg. mIoU	SrcSep. SDR
✗	36.82	45.29	8.79	31.57	38.73	55.32	21.38	3.52
✓	45.38	49.72	15.56	34.83	41.96	60.82	24.62	4.63

ison results in Table 9. Our model showcased the best performance in both SDR and SAR across all datasets. For example, we achieved an SDR score of 10.75 on the MUSIC dataset, significantly higher than other methods, reflecting superior separation quality. Meanwhile, our model also reached an SDR score of 8.23 and 5.06 on VGGSound-Instruments and VGGSound-Music, respectively. These results underscore the effectiveness of *ACAV-1M* for audio-visual source separation.

4.3 EXPERIMENTAL ANALYSIS

In this section, we provide a detailed analysis of our experimental results to demonstrate the benefits of our data curation pipeline, the impact of the quality measure criteria on model performance, and the scaling behavior of the *ACAV-1M* dataset.

Benefit of data curation pipeline. To demonstrate the efficacy of our data curation pipeline, we conducted comparative experiments using a random subset of VGGSound with equivalent size and a clean subset of VGGSound with our data curation pipeline. The experimental results are reported in Table 10. The results clearly indicate that improvements are attributed to our data curation process, which includes precise alignment of audio and visual data and careful annotation. The clean subset of VGGSound with our data curation pipeline achieves superior performance and demonstrates the value of high-quality, well-aligned data in training more effective multimodal models.

Ablation on alignment filtering. We analyzed the impact of different quality measure criteria on model performance, as shown in Table 11. For language alignment, we experimented with using class names directly instead of our detailed, long captions for annotations. This change resulted in a noticeable degradation in performance, emphasizing the importance of rich, descriptive captions in providing contextual cues that enhance model understanding and performance. Regarding temporal alignment, we varied the alignment accuracy of audio and visual data during the dataset curation process, testing alignment accuracies of 100%, 70%, and 50%. Our experiments show that models trained with 100% alignment accuracy consistently outperform those trained with lower accuracies, underscoring the critical role of precise synchronization in audio-visual learning.

Ablation on similarity threshold. The similarity threshold of the *ACAV-1M* filtering across different audio-visual benchmarks reveals important insights into the model’s performance sensitivity to the alignment accuracy between audio and visual modalities. The ablation results are reported in Table 12. At a 75% similarity threshold, the results demonstrate moderate performance across all metrics, with classification accuracy and source localization precision peaking at 46.23% and 49.52%, respectively. This higher threshold suggests that stringent filtering criteria may exclude valuable, albeit less perfectly aligned, audio-visual pairs which could contribute useful information to the tasks. Lowering the threshold to 50% yields the best overall performance with notable improvements in several key areas: classification accuracy increases to 47.85%, and retrieval accuracy reaches its peak at 16.78%. However, further reduction of the threshold to 25% results in a general decrease in performance across most metrics, with classification accuracy dropping to 43.72%, and a

Table 11: Ablation results on quality measure criteria across different audio-visual benchmarks.

Quality Measure	Cls. Acc (%)	SrcLoc. Prec	Retrieval Acc (%)	SADialog. BLEU	VidPars. F-score (%)	QA Acc (%)	Seg. mIoU	SrcSep. SDR
–	33.25	40.67	4.86	29.78	36.93	50.19	15.86	1.87
Instance Align	39.56	43.95	6.17	31.23	38.21	53.27	18.37	2.38
+ Temporal Align	42.68	46.53	12.65	33.15	40.16	55.32	20.65	3.29
+ Language Align	47.85	50.27	16.78	35.34	42.65	60.26	24.83	4.35
threshold=70%	45.76	49.12	15.69	34.87	42.36	59.75	24.02	4.13
threshold=50%	44.65	48.35	14.73	34.06	41.89	58.23	22.96	3.67

Table 12: Ablation results on the similarity threshold of our ACAV-1M filtering across different audio-visual benchmarks.

Similarity Threshold	Cls. Acc (%)	SrcLoc. Prec	Retrieval Acc (%)	SADialog. BLEU	VidPars. F-score (%)	QA Acc (%)	Seg. mIoU	SrcSep. SDR
75%	46.23	49.52	15.85	34.96	42.56	59.83	24.37	4.19
50%	47.85	50.27	16.78	35.34	42.65	60.26	24.83	4.35
25%	43.72	47.67	14.89	33.67	40.23	57.08	21.25	2.58

Table 13: Ablation results on the scaling property of our ACAV-1M dataset across different audio-visual benchmarks.

Data Scale	Cls. Acc (%)	SrcLoc. Prec	Retrieval Acc (%)	SADialog. BLEU	VidPars. F-score (%)	QA Acc (%)	Seg. mIoU	SrcSep. SDR
10K	30.63	37.82	4.23	28.56	33.25	45.63	11.56	1.25
52K	47.85	50.27	16.78	35.34	42.65	60.26	24.83	4.35
100K	49.27	51.85	17.63	36.08	43.72	62.23	25.98	5.12
199K	51.73	53.62	19.35	38.15	46.58	65.73	28.75	6.28
1M	61.35	60.75	26.57	43.27	55.35	73.25	36.39	10.75

significant decrease in segmentation performance (mIoU) to 21.25% and source separation (SDR) to 2.58. This suggests that too low a threshold includes too many poorly aligned pairs, which confuses the model and degrades the overall performance.

Scaling property of our ACAV-1M dataset. To understand the scalability of our dataset, we trained models on progressively larger subsets of *ACAV-1M*, specifically 10K, 52K, 100K, 199K, and 1M samples. The experimental results across different downstream tasks are reported in Table 13. Our findings reveal a positive correlation between the size of the dataset and the performance. As the size increases, we observe improved accuracy and robustness across all tasks, indicating that *ACAV-1M* not only supports effective training at smaller scales but also benefits significantly from scaling up.

5 CONCLUSION

In this work, we introduce *ACAV-1M*, a novel large-scale dataset with one million samples that are curated to bridge the gap between audio and visual data. Furthermore we propose a comprehensive audio-visual benchmark that supports a wide array of audio-visual tasks, from classification and segmentation to retrieval and scene-aware dialog, each benefiting from the dataset’s rich annotations and precise audio-visual alignment. We demonstrate the superior performance of models trained on *ACAV-1M* compared to existing methods. Our experiments also explore the scaling behavior of the dataset, showing significant improvements in model performance as data volume increases, thus confirming the dataset’s scalability.

Limitations and broader impact. While our dataset covers a broad range of audio and visual contexts, there are still rare scenarios that are underrepresented or absent, which could affect the generalizability of the trained models to all real-world applications.

ACAV-1M has the potential to make a profound impact for audio-visual learning. The insights gained from models trained on *ACAV-1M* can enhance multimedia applications, improve accessibility features, and foster the development of intuitive and interactive systems. However, it is essential to be aware of ethical considerations and potential biases in training data, which could amplify disparities if not carefully managed.

ETHICS STATEMENT

In accordance with the ICLR Code of Ethics, our research strictly utilizes datasets that are publicly available and have been released for academic use. We recognize the implications of deploying machine learning models in the real world, particularly concerning the potential for unintended consequences. Therefore, we stress the importance of using our findings and methodologies responsibly. Our experiments are designed to foster advancements in audio-visual processing technologies while ensuring that these technologies are developed in an ethical manner. We are open to discussions regarding the ethical considerations of our work and actively seek feedback to refine our approach in line with best practices.

REPRODUCIBILITY STATEMENT

To ensure the integrity and reproducibility of our research, we have meticulously documented our algorithms, experimental design, and implementation details in Section 3, Section 4 and Appendix F of our submission. Post-publication, we are committed to making the codebase publicly available, which encompasses all pertinent scripts and models used in our experiments. This open-source approach is intended to facilitate validation of our results by the broader scientific community and to support future research endeavors that build upon our work.

REFERENCES

- Huda AlAmri, Vincent Cartillier, Abhishek Das, Jue Wang, Stefan Lee, Peter Anderson, Irfan Essa, Devi Parikh, Dhruv Batra, Anoop Cherian, Tim K. Marks, and Chiori Hori. Audio visual scene-aware dialog. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7550–7559, 2019. [2](#)
- Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 609–617, 2017. [1](#), [2](#)
- Relja Arandjelović and Andrew Zisserman. Look, listen and learn. *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 609–617, 2017. [2](#)
- Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 435–451, 2018. [6](#), [7](#), [8](#)
- Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2016. [1](#), [2](#)
- Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6299–6308, 2017. [24](#)
- Changan Chen, Unnat Jain, Carl Schissler, S. V. A. Garí, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman. Soundspaces: Audio-visual navigation in 3d environments. In *Proceedings of European Conference on Computer Vision (ECCV)*, pp. 17–36, 2020a. [2](#)
- Changan Chen, Sagnik Majumder, Al-Halah Ziad, Ruohan Gao, Santhosh Kumar Ramakrishnan, and Kristen Grauman. Learning to set waypoints for audio-visual navigation. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021a. [2](#)
- Changan Chen, Carl Schissler, Sanchit Garg, Philip Kobernik, Alexander Clegg, Paul Calamia, Dhruv Batra, Philip W Robinson, and Kristen Grauman. Soundspaces 2.0: A simulation platform for visual-acoustic learning. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2022. [2](#)
- Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 721–725. IEEE, 2020b. [2](#), [3](#), [5](#), [25](#)

- 594 Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zis-
595 serman. Localizing visual sounds the hard way. In *Proceedings of the IEEE/CVF Conference on*
596 *Computer Vision and Pattern Recognition (CVPR)*, pp. 16867–16876, 2021b. [6](#), [7](#), [8](#)
597
- 598 Jia Deng, Wei Dong, Richard Socher, Li-Jia. Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale
599 Hierarchical Image Database. In *Proceedings of IEEE/CVF Conference on Computer Vision and*
600 *Pattern Recognition (CVPR)*, pp. 248–255, 2009. [24](#)
- 601 Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: an audio captioning dataset.
602 *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing*
603 *(ICASSP)*, pp. 736–740, 2019. [2](#)
604
- 605 Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T
606 Freeman, and Michael Rubinstein. Looking to listen at the cocktail party: A speaker-independent
607 audio-visual model for speech separation. *arXiv preprint arXiv:1804.03619*, 2018. [2](#)
- 608 Chuang Gan, Deng Huang, Hang Zhao, Joshua B. Tenenbaum, and Antonio Torralba. Music gesture
609 for visual sound separation. In *IEEE/CVF Conference on Computer Vision and Pattern Recogni-*
610 *tion (CVPR)*, pp. 10478–10487, 2020. [2](#)
611
- 612 Ruohan Gao and Kristen Grauman. 2.5d visual sound. In *Proceedings of the IEEE/CVF Conference*
613 *on Computer Vision and Pattern Recognition (CVPR)*, pp. 324–333, 2019. [2](#)
- 614 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,
615 Hal Dauméé III, and Kate Crawford. Datasheets for Datasets. *arXiv preprint arXiv:1803.09010*,
616 2018. [17](#)
617
- 618 Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing
619 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for
620 audio events. In *Proceedings of 2017 IEEE International Conference on Acoustics, Speech and*
621 *Signal Processing (ICASSP)*, pp. 776–780, 2017. [2](#), [3](#), [24](#), [25](#), [33](#)
- 622 Mariana-Iuliana Georgescu, Eduardo Fonseca, Radu Tudor Ionescu, Mario Lucic, Cordelia Schmid,
623 and Anurag Arnab. Audiovisual masked autoencoders. In *Proceedings of the IEEE/CVF Interna-*
624 *tional Conference on Computer Vision (ICCV)*, pp. 16144–16154, October 2023. [5](#), [6](#)
- 625 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
626 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of*
627 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. [1](#), [3](#), [4](#)
628
- 629 Yuan Gong, Andrew Rouditchenko, Alexander H. Liu, David Harwath, Leonid Karlinsky, Hilde
630 Kuehne, and James R. Glass. Contrastive audio-visual masked autoencoder. In *Proceedings of*
631 *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. [5](#), [6](#)
- 632 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
633 nition. In *Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition*
634 *(CVPR)*, pp. 770–778, 2016. [24](#)
635
- 636 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked
637 autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. [5](#), [6](#), [24](#)
- 638 John Hershey and Michael Casey. Audio-visual sound separation via hidden markov models. *Ad-*
639 *vances in Neural Information Processing Systems*, 14, 2001. [2](#)
640
- 641 Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing
642 Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss,
643 and Kevin W. Wilson. CNN architectures for large-scale audio classification. In *Proceedings of*
644 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017. [24](#)
- 645 Chiori Hori, Huda Alamri, Jue Wang, Gordon Winchern, Takaaki Hori, Anoop Cherian, Tim K
646 Marks, Vincent Cartillier, Raphael Gontijo Lopes, Abhishek Das, et al. End-to-end audio
647 visual scene-aware dialog using multimodal attention-based video features. *arXiv preprint*
arXiv:1806.08409, 2018. [6](#), [7](#)

- 648 Di Hu, Feiping Nie, and Xuelong Li. Deep multimodal clustering for unsupervised audiovisual
649 learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
650 *(CVPR)*, pp. 9248–9257, 2019. 2, 6, 7, 8
- 651
- 652 Di Hu, Rui Qian, Minyue Jiang, Xiao Tan, Shilei Wen, Errui Ding, Weiyao Lin, and Dejing Dou.
653 Discriminative sounding objects localization via self-supervised audiovisual matching. In *Pro-*
654 *ceedings of Advances in Neural Information Processing Systems (NeurIPS)*, pp. 10077–10087,
655 2020. 6, 7, 8
- 656 Xixi Hu, Ziyang Chen, and Andrew Owens. Mix and localize: Localizing sound sources in mix-
657 tures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
658 *(CVPR)*, pp. 10483–10492, 2022. 5, 6, 7, 8, 24, 25
- 659 Po-Sen Huang, Scott Deeann Chen, Paris Smaragdis, and Mark Hasegawa-Johnson. Singing-voice
660 separation from monaural recordings using robust principal component analysis. In *IEEE Inter-*
661 *national Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 57–60, 2012. 8,
662 9
- 663
- 664 Po-Yao Huang, Vasu Sharma, Hu Xu, Chaitanya Ryali, Haoqi Fan, Yanghao Li, Shang-Wen Li,
665 Gargi Ghosh, Jitendra Malik, and Christoph Feichtenhofer. Mavil: Masked audio-video learners.
666 *arXiv preprint arXiv:2212.08071*, 2022a. 5, 6
- 667 Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian
668 Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *Proceedings of Ad-*
669 *vances In Neural Information Processing Systems (NeurIPS)*, 2022b. 5, 6, 24
- 670
- 671 Vladimir Iashin and Esa Rahtu. A better use of audio-visual cues: Dense video captioning with
672 bi-modal transformer. In *British Machine Vision Conference (BMVC)*, 2020. 6, 7
- 673 Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating
674 captions for audios in the wild. In *NAACL-HLT*, 2019. 2
- 675
- 676 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
677 *arXiv:1412.6980*, 2014. 6, 24, 25
- 678 Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models
679 from self-supervised synchronization. In *Proceedings of Advances in Neural Information Pro-*
680 *cessing Systems (NeurIPS)*, 2018. 1, 2
- 681
- 682 Sangho Lee, Jiwan Chung, Youngjae Yu, Gunhee Kim, Thomas Breuel, Gal Chechik, and Yale
683 Song. Acav100m: Automatic curation of large-scale datasets for audio-visual video representa-
684 tion learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*,
685 2021. 2, 3, 18, 19, 20
- 686 Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji rong Wen, and Di Hu. Learning to an-
687 swer questions in dynamic audio-visual scenarios. In *Proceedings of IEEE/CVF Conference on*
688 *Computer Vision and Pattern Recognition (CVPR)*, pp. 19086–19096, 2022. 2, 5, 8, 26
- 689
- 690 Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning
691 united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*,
692 2023. 1, 3, 4, 26
- 693 Yan-Bo Lin and Yu-Chiang Frank Wang. Audiovisual transformer with instance attention for audio-
694 visual event localization. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*,
695 2020. 2
- 696
- 697 Yan-Bo Lin, Yu-Jhe Li, and Yu-Chiang Frank Wang. Dual-modality seq2seq network for audio-
698 visual event localization. In *IEEE International Conference on Acoustics, Speech and Signal*
699 *Processing (ICASSP)*, pp. 2002–2006, 2019. 2, 7, 8
- 700 Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, and Ming-Hsuan Yang. Exploring cross-
701 video and cross-modality signals for weakly-supervised audio-visual video parsing. In *Proce-*
edings of Advances in Neural Information Processing Systems (NeurIPS), 2021. 2

- 702 Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V
703 Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective
704 instruction tuning. *arXiv preprint arXiv:2301.13688*, 2023. [1](#), [2](#)
705
- 706 Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D. Plumb-
707 ley, Yuexian Zou, and Wenwu Wang. WavCaps: A ChatGPT-assisted weakly-labelled audio
708 captioning dataset for audio-language multimodal research. *IEEE/ACM Transactions on Audio,*
709 *Speech, and Language Processing*, pp. 1–15, 2024. [3](#), [4](#), [26](#), [34](#)
- 710 Shentong Mo and Pedro Morgado. Localizing visual sounds the easy way. In *Proceedings of*
711 *European Conference on Computer Vision (ECCV)*, pp. 218–234, 2022a. [6](#), [7](#), [8](#), [24](#)
712
- 713 Shentong Mo and Pedro Morgado. A closer look at weakly-supervised audio-visual source localiza-
714 tion. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022b.
715 [5](#), [6](#), [7](#), [8](#), [24](#)
- 716 Shentong Mo and Pedro Morgado. A unified audio-visual learning framework for localization,
717 separation, and recognition. In *Proceedings of the International Conference on Machine Learning*
718 *(ICML)*, 2023. [5](#), [8](#), [9](#), [25](#)
719
- 720 Shentong Mo and Yapeng Tian. Multi-modal grouping network for weakly-supervised audio-visual
721 video parsing. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*,
722 2022. [2](#), [7](#), [8](#)
- 723 Pedro Morgado, Nuno Vasconcelos, Timothy Langlois, and Oliver Wang. Self-supervised generation
724 of spatial audio for 360 video. In *Advances in Neural Information Processing Systems (NeurIPS)*,
725 2018. [2](#)
726
- 727 Pedro Morgado, Yi Li, and Nuno Vasconcelos. Learning representations from audio-visual spatial
728 alignment. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, pp.
729 4733–4744, 2020. [2](#)
- 730 Pedro Morgado, Ishan Misra, and Nuno Vasconcelos. Robust audio-visual instance discrimina-
731 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
732 *(CVPR)*, pp. 12934–12945, 2021a. [2](#)
733
- 734 Pedro Morgado, Nuno Vasconcelos, and Ishan Misra. Audio-visual instance discrimination with
735 cross-modal agreement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
736 *Pattern Recognition (CVPR)*, pp. 12475–12486, June 2021b. [2](#)
- 737 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. [1](#), [3](#), [4](#), [26](#)
738
- 739 Andrew Owens, Jiajun Wu, Josh H. McDermott, William T. Freeman, and Antonio Torralba. Am-
740 bient sound provides supervision for visual learning. In *Proceedings of the European Conference*
741 *on Computer Vision (ECCV)*, pp. 801–816, 2016. [1](#), [2](#)
742
- 743 Rui Qian, Di Hu, Heinrich Dinkel, Mengyue Wu, Ning Xu, and Weiyao Lin. Multiple sound sources
744 localization from coarse to fine. In *Proceedings of European Conference on Computer Vision*
745 *(ECCV)*, pp. 292–308, 2020. [6](#), [7](#), [8](#)
- 746 Andrew Rouditchenko, Angie Boggust, David Harwath, Brian Chen, Dhiraj Joshi, Samuel Thomas,
747 Kartik Audhkhasi, Hilde Kuehne, Rameswar Panda, Rogerio Feris, et al. Avlnet: Learning audio-
748 visual language representations from instructional videos. *arXiv preprint arXiv:2006.09199*,
749 2020. [6](#), [7](#)
- 750 Idan Schwartz, Alexander G. Schwing, and Tamir Hazan. A simple baseline for audio-visual scene-
751 aware dialog. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
752 *Recognition (CVPR)*, 2019. [8](#)
753
- 754 Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to local-
755 ize sound source in visual scenes. In *Proceedings of the IEEE Conference on Computer Vision*
and Pattern Recognition (CVPR), pp. 4358–4366, 2018. [1](#), [2](#), [4](#), [6](#), [7](#), [8](#), [25](#)

- 756 SAnkit Shah, Shijie Geng, Peng Gao, Anoop Cherian, Takaaki Hori, Tim K. Marks, Jonathan Le
757 Roux, and Chiori Hori. Audio-visual scene-aware dialog and reasoning using audio-visual trans-
758 formers with joint student-teacher learning. In *Proceedings of IEEE International Conference on*
759 *Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7732–7736. IEEE, 2022. 5, 6, 7
- 760 Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and Policy Considerations for
761 Deep Learning in NLP. *arXiv preprint arXiv:1906.02243*, 2019. 20
- 762 Zineng Tang, Jaemin Cho, Yixin Nie, and Mohit Bansal. Tvlv: Textless vision-language transformer.
763 In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 6, 7
- 764 Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization
765 in unconstrained videos. In *Proceedings of European Conference on Computer Vision (ECCV)*,
766 2018. 2, 7, 8
- 767 Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-supervised
768 audio-visual video parsing. In *Proceedings of European Conference on Computer Vision (ECCV)*,
769 pp. 436–454, 2020. 2, 5, 7, 8, 24, 26
- 770 Yapeng Tian, Di Hu, and Chenliang Xu. Cyclic co-learning of sounding object visual grounding and
771 sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
772 *Recognition (CVPR)*, pp. 2745–2754, 2021. 8, 9
- 773 Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer
774 look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE Confer-*
775 *ence on Computer Vision and Pattern Recognition (CVPR)*, pp. 6450–6459, 2018. 24
- 776 Tuomas Virtanen. Monaural sound source separation by nonnegative matrix factorization with tem-
777 poral continuity and sparseness criteria. *IEEE Transactions on Audio, Speech, and Language*
778 *Processing*, 15(3):1066–1074, 2007. 8, 9
- 779 Zehan Wang, Ziang Zhang, Xize Cheng, Rongjie Huang, Luping Liu, Zhenhui Ye, Haifeng Huang,
780 Yang Zhao, Tao Jin, Peng Gao, and Zhou Zhao. Freebind: Free lunch in unified multimodal space
781 via knowledge fusion. In *Proceedings of the International Conference on Machine Learning*
782 *(ICML)*, 2024a. 34
- 783 Zehan Wang, Ziang Zhang, Hang Zhang, Luping Liu, Rongjie Huang, Xize Cheng, Hengshuang
784 Zhao, and Zhou Zhao. Omnibind: Large-scale omni multimodal representation via binding
785 spaces. *arXiv preprint arXiv:2407.11895*, 2024b. 34
- 786 Yu Wu and Yi Yang. Exploring heterogeneous clues for weakly-supervised audio-visual video pars-
787 ing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
788 pp. 1326–1335, 2021. 2
- 789 Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention matching for audio-visual event local-
790 ization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp.
791 6291–6299, 2019. 2
- 792 Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridg-
793 ing video and language. In *Proceedings of IEEE Conference on Computer Vision and Pattern*
794 *Recognition (CVPR)*, pp. 5288–5296, 2016. 5, 25
- 795 Xudong Xu, Bo Dai, and Dahua Lin. Recursive visual sound separation using minus-plus net. In
796 *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 8, 9
- 797 Heeseung Yun, Youngjae Yu, Wonsuk Yang, Kangil Lee, and Gunhee Kim. Pano-avqa: Grounded
798 audio-visual question answering on 360° videos. In *Proceedings of the IEEE International Con-*
799 *ference on Computer Vision (ICCV)*, 2021. 8
- 800 Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio
801 Torralba. The sound of pixels. In *Proceedings of the European Conference on Computer Vision*
802 *(ECCV)*, pp. 570–586, 2018. 2, 5, 6, 8, 9, 25
- 803 Hang Zhao, Chuang Gan, Wei-Chiu Ma, and Antonio Torralba. The sound of motions. In *Proced-*
804 *ings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 1735–1744,
805 2019. 2

810 Jinxing Zhou, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birchfield, Dan Guo,
811 Lingpeng Kong, Meng Wang, and Yiran Zhong. Audio-visual segmentation. In *Proceedings of*
812 *European Conference on Computer Vision (ECCV)*, 2022. [2](#), [5](#), [6](#), [8](#), [25](#)
813
814 Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Young-
815 jae Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal C4: An open, billion-
816 scale corpus of images interleaved with text. *arXiv preprint arXiv:2304.06939*, 2023. [1](#), [2](#)
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
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864 APPENDIX
865

866 In this appendix, we provide the following material:
867

- 868 • dataset documentation and intended uses in Section A,
- 869 • dataset website in Section B,
- 870 • croissant metadata in Section C,
- 871 • author statement in Section D,
- 872 • [licensing overview and details in Section E](#),
- 873 • addition implementation and datasets details in Section F,
- 874 • [additional data construction details and quality analysis in Section G](#),
- 875 • [additional experimental analyses in Section H](#).

876
877
878 A DATASET DOCUMENTATION & INTENDED USES
879

880 The *ACAV-1M* dataset is designed to facilitate research in audio-visual representation learning. It
881 includes synchronized audio and visual data curated from various sources to ensure a diverse and
882 comprehensive collection for training and evaluating machine learning models. The dataset doc-
883 umentation follows the *datasheets for datasets* framework, providing detailed information on the
884 dataset’s composition, collection process, and intended uses.

885 Composition: The dataset consists of 1 million audio-visual pairs, including video clips with corre-
886 sponding audio tracks from various domains such as user-generated content.
887

888 Collection Process: Data was collected using automated scripts and manual curation to ensure qual-
889 ity and relevance. Metadata includes source URLs, timestamps, and content descriptions.

890 Intended Uses: The dataset is intended for developing and benchmarking models in audio-visual
891 representation learning, including tasks like video classification, audio-visual synchronization, and
892 cross-modal retrieval.

893 Ethical Considerations: We ensured that the dataset adheres to ethical guidelines, including the
894 exclusion of sensitive or inappropriate content and respect for copyright and privacy concerns.

895 **This document is based on *Datasheets for Datasets* by Gebru *et al.* (Gebru *et al.*, 2018).**
896

897
898 **MOTIVATION**
899

900 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a
901 specific gap that needed to be filled? Please provide a description.

902 The *ACAV-1M* was created to fill a significant gap in multimodal learning where audio and visual
903 data are integrated systematically. It aims to enhance robust models that leverage both modalities
904 for improved understanding and interaction, designed specifically for tasks like audio-visual
905 classification, localization, retrieval, and segmentation.

906 **Who created this dataset (e.g., which team, research group) and on behalf of which entity**
907 **(e.g., company, institution, organization)?**

908 The dataset was created by a collaborative effort involving researchers from various academic
909 institutions specializing in machine learning and computer vision, under the coordination of a
910 leading university’s computer science department.
911

912 **What support was needed to make this dataset?** (e.g. who funded the creation of the dataset? If
913 there is an associated grant, provide the name of the grantor and the grant name and number, or if it
914 was supported by a company or government agency, give those details.)

915 No. The creation of the *ACAV-1M* was not supported by any grants from several research funding
916 agencies. However, the dataset development received technical support and infrastructure from the
917 host university.

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

Any other comments?

No.

COMPOSITION

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

The instances in the *ACAV-1M* represent synchronized audio-visual clips from diverse settings, including music performances, public speeches, and everyday activities, ensuring a wide range of scenarios for robust multimodal learning.

How many instances are there in total (of each type, if appropriate)?

The dataset comprises approximately 100,000 video clips, each paired with corresponding audio tracks that have been meticulously synchronized and annotated.

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

The dataset is a curation subset of the original ACAV (Lee et al., 2021) dataset with 100 million samples.

What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or features? In either case, please provide a description.

Each instance consists of “Raw” video and audio data. Additional metadata include synchronization points, annotations for source localization, and labels for classification and segmentation tasks.

Is there a label or target associated with each instance? If so, please provide a description.

Yes, each instance includes captions associated with the video and audio. For various tasks, we include labels for each instance like classification (audio-visual context), segmentation masks, and localization coordinates.

Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

No.

Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network links)? If so, please describe how these relationships are made explicit.

No.

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

No.

Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.

No.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? If it links to or relies on external resources, a) are there

972 guarantees that they will exist, and remain constant, over time; b) are there official archival versions
973 of the complete dataset (i.e., including the external resources as they existed at the time the dataset
974 was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external
975 resources that might apply to a future user? Please provide descriptions of all external resources
976 and any restrictions associated with them, as well as links or other access points, as appropriate.

977 Yes. The dataset is a curation subset of the original ACAV (Lee et al., 2021) dataset with 100
978 million samples.

979
980 **Does the dataset contain data that might be considered confidential (e.g., data that is**
981 **protected by legal privilege or by doctor-patient confidentiality, data that includes the content**
982 **of individuals' non-public communications)?** If so, please provide a description.
983 No.

984
985 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-**
986 **ing, or might otherwise cause anxiety?** If so, please describe why.
987 No.

988
989 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.
990 No.

991
992 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how
993 these subpopulations are identified and provide a description of their respective distributions within
994 the dataset.
995 No.

996
997 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or**
998 **indirectly (i.e., in combination with other data) from the dataset?** If so, please describe how.
999 No.

1000
1001 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that**
1002 **reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or**
1003 **union memberships, or locations; financial or health data; biometric or genetic data; forms of**
1004 **government identification, such as social security numbers; criminal history)?** If so, please
1005 provide a description.
1006 No.

1007
1008 **Any other comments?**
1009 No.

1010
1011 **COLLECTION**

1012
1013
1014 **How was the data associated with each instance acquired?** Was the data directly observable
1015 (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly in-
1016 ferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or language)?
1017 If data was reported by subjects or indirectly inferred/derived from other data, was the data
1018 validated/verified? If so, please describe how.

1019 The data was acquired through a combination of public domain resources and contributions from
1020 collaborating institutions, where scenarios were staged and recorded under controlled conditions to
1021 ensure quality and diversity.

1022
1023 **Over what timeframe was the data collected?** Does this timeframe match the creation timeframe
1024 of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please
1025 describe the timeframe in which the data associated with the instances was created. Finally, list
when the dataset was first published.

Data collection spanned over half one year, culminating in the dataset's release in 2024. The

1026 temporal alignment of collection and creation ensured the relevance and recency of the data.
1027

1028 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus**
1029 **or sensor, manual human curation, software program, software API)?** How were these
1030 mechanisms or procedures validated?

1031 We have alignment filtering mechanisms to curate our dataset from the original ACAV (Lee et al.,
1032 2021) dataset.
1033

1034 **What was the resource cost of collecting the data?** (e.g. what were the required computational
1035 resources, and the associated financial costs, and energy consumption - estimate the carbon
1036 footprint. See Strubell *et al.* Strubell et al. (2019) for approaches in this area.)

1037 We use four A100 GPUs to curate data and train our models.
1038

1039 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
1040 **probabilistic with specific sampling probabilities)?**

1041 We used alignment filtering mechanisms.
1042

1043 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
1044 **and how were they compensated (e.g., how much were crowdworkers paid)?**

1045 Authors are involved in the data curation process.
1046

1047 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so,
1048 please provide a description of these review processes, including the outcomes, as well as a link or
1049 other access point to any supporting documentation.

1050 No.
1051

1052 **Does the dataset relate to people?** If not, you may skip the remainder of the questions in this
1053 section.

1054 No.
1055

1056 **Did you collect the data from the individuals in question directly, or obtain it via third parties**
1057 **or other sources (e.g., websites)?**

1058 No.
1059

1060 **Were the individuals in question notified about the data collection?** If so, please describe (or
1061 show with screenshots or other information) how notice was provided, and provide a link or other
1062 access point to, or otherwise reproduce, the exact language of the notification itself.

1063 No.
1064

1065 **Did the individuals in question consent to the collection and use of their data?** If so, please
1066 describe (or show with screenshots or other information) how consent was requested and provided,
1067 and provide a link or other access point to, or otherwise reproduce, the exact language to which the
1068 individuals consented.

1069 No.
1070

1071 **If consent was obtained, were the consenting individuals provided with a mechanism to**
1072 **revoke their consent in the future or for certain uses?** If so, please provide a description, as
1073 well as a link or other access point to the mechanism (if appropriate)

1074 No.
1075

1076 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data**
1077 **protection impact analysis) been conducted?** If so, please provide a description of this analysis,
1078 including the outcomes, as well as a link or other access point to any supporting documentation.
1079

1080 No.

1081

1082

Any other comments?

1083

No.

1084

1085

1086

PREPROCESSING / CLEANING / LABELING

1087

1088

1089

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)? If so, please provide a description. If not, you may skip the remainder of the questions in this section.

1092

Yes. We use multimodal LLM to .

1093

1094

1095

Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.

1096

No.

1097

1098

1099

Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.

1100

No.

1101

1102

1103

Any other comments?

1104

No.

1105

1106

USES

1107

1108

1109

Has the dataset been used for any tasks already? If so, please provide a description.

1110

Yes, *ACAV-1M* has been employed in several benchmarking tasks within the research group, including preliminary studies on audio-visual perception tasks.

1111

1112

1113

Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.

1114

No.

1115

1116

1117

What (other) tasks could the dataset be used for?

1118

Beyond the current uses, the dataset holds potential for tasks in automated content generation, assistive technologies, and advanced surveillance systems.

1119

1120

1121

Is there anything about the composition of the dataset or the way it was collected and pre-processed/cleaned/labeled that might impact future uses? For example, is there anything that a future user 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 undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is there anything a future user could do to mitigate these undesirable harms?

1122

No.

1123

1124

1125

Are there tasks for which the dataset should not be used? If so, please provide a description.

1129

No.

1130

1131

Any other comments?

1132

No.

1133

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

DISTRIBUTION

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? If so, please provide a description.
No.

How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?
The dataset is available via a website page and can be accessed through the dataset page, which ensures controlled and ethical usage aligned with academic standards.

When will the dataset be distributed?
The dataset will be available upon publication.

Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.
No.

Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.
No.

Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.
No.

Any other comments?
No.

MAINTENANCE

Who is supporting/hosting/maintaining the dataset?
The dataset is maintained by the authors, with plans for ongoing updates and expansions based on community feedback and technological advancements.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)?
The owner of the dataset can contacted by email.

Is there an erratum? If so, please provide a link or other access point.
No.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)? If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?
Yes, the dataset is scheduled for regular reviews and updates to address any errors, introduce new instances, and phase out obsolete data, with all changes communicated through the dataset's official repository.

1188 **If the dataset relates to people, are there applicable limits on the retention of the data**
 1189 **associated with the instances (e.g., were individuals in question told that their data would be**
 1190 **retained for a fixed period of time and then deleted)?** If so, please describe these limits and
 1191 explain how they will be enforced.
 1192 No.

1193
 1194 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please
 1195 describe how. If not, please describe how its obsolescence will be communicated to users.
 1196 Yes. It will be maintained on the dataset website.

1197
 1198 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
 1199 **them to do so?** If so, please provide a description. Will these contributions be validated/verified?
 1200 If so, please describe how. If not, why not? Is there a process for communicating/distributing these
 1201 contributions to other users? If so, please provide a description.
 1202 Yes. We will open the opportunity for other researchers to augment the dataset for additional
 1203 benchmarks.

1204
 1205 **Any other comments?**
 1206 No.

1207 1208 1209 B DATASET WEBSITE

1210
 1211 The dataset and its documentation can be accessed at the following URL:

1212 <https://acav1m.github.io>
 1213

1214 This website provides an overview of the dataset, download links, and additional resources such
 1215 as example code, tutorials, and a forum for community discussions. Users can explore the dataset
 1216 through an interactive interface, which includes search and filter options to facilitate easy access to
 1217 specific subsets of the data.

1218 1219 C CROISSANT METADATA

1220
 1221 The Croissant metadata for the *ACAV-1M* dataset is available at:

1222 <https://acav1m.github.io>
 1223

1224 This metadata record documents the dataset’s structure, including descriptions of the files, their
 1225 formats, and the fields within each record. The metadata adheres to the Croissant format, ensuring
 1226 interoperability and ease of use with ML tools and platforms. [We have structured metadata](#)
 1227 [for *ACAV-1M* including video captions, timestamps, and similarity scores.](#) [Table 14 provides an](#)
 1228 [overview of the metadata fields available to researchers.](#)

1229
 1230 Table 14: Metadata Fields in *ACAV-1M*.

Field	Description
Video ID	Unique identifier for each video
Timestamps	Start and end times for clips
Captions	Aggregated audio-visual descriptions
Similarity Scores	Text-audio, text-visual, audio-visual alignments
Categories	High-level category labels

1237 1238 1239 D AUTHOR STATEMENT

1240
 1241 We, the authors of the *ACAV-1M* dataset, bear full responsibility for any violations of rights and confirm that all data included in the dataset complies with the relevant licenses and ethical guidelines.

The dataset is released under the Creative Commons Attribution 4.0 International License (CC BY 4.0), which allows for sharing, adaptation, and use of the data with appropriate credit given to the original authors.

E LICENSING OVERVIEW

The dataset is licensed under the Creative Commons Attribution 4.0 International License. To ensure ethical compliance, we sourced videos under licenses permitting academic use. Table 15 outlines the proportion of data from each source and its associated licensing terms.

Table 15: Licensing Details for Dataset Sources.

Source	License Type	Proportion (%)
YouTube	Creative Commons	60
Existing AV Datasets	Academic Research Agreements	30
Public Repositories	Open-Source Licenses	10

F IMPLEMENTATION & DATASET DETAILS

In this section, we provide more implementation and dataset details.

Audio-visual classification. For linear probing, we follow the prior work (He et al., 2021; Huang et al., 2022b) and extract frozen audio-visual representations from our *ACAV-1M* pre-trained audio-visual masked autoencoder. Then we attach a linear layer as a head to the frozen features for training with the audio-visual classes. During training, we only fine-tune the linear head to evaluate the quality of pre-trained features. The models are trained for 50 epochs using the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $1e-4$ and a batch size of 128. For fine-tuning, we use the same optimizer and batch size settings, but all parameters are learnable.

Audio-Visual Source Localization. For sound source localization, we train all baselines (Mo & Morgado, 2022a;b; Hu et al., 2022) using the same backbone (*i.e.*, ViT-Base) for audio/visual encoder with different proposed objectives in their original papers. The final localization map is generated through bilinear interpolation of the similarity map between audio/visual features from the last self-attention layer. The models are trained for 30 epochs using the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $1e-4$ and a batch size of 128.

Audio-Visual Retrieval. The retrieval task processes video frames sampled at 8 fps and utilizes combined low-level visual features from ResNet-152 (He et al., 2016) and 3D ResNet models (Tran et al., 2018), both pre-trained on respective large-scale datasets. Audio features are extracted using VGGish (Hershey et al., 2017), pre-trained on AudioSet (Gemmeke et al., 2017). The complete model, integrating these features, is trained using Adam to optimize retrieval effectiveness across 40 epochs.

Audio-Visual Video Parsing. Following the data pre-processing in previous work (Tian et al., 2020), we sample video frames at 8 fps from the 10-second videos with 10 non-overlapping snippets of 1 second. For low-level visual features, we concatenate 2D and 3D visual features extracted by ResNet-152 (He et al., 2016) pre-trained on ImageNet (Deng et al., 2009) and 3D ResNet (Tran et al., 2018) pre-trained on Kinetics-400 (Carreira & Zisserman, 2017). We utilize VGGish (Hershey et al., 2017) pre-trained on AudioSet (Gemmeke et al., 2017) to extract the audio features. The model is trained with Adam (Kingma & Ba, 2014) optimizer with $\beta_1=0.9$, $\beta_2=0.999$ and with an initial learning rate of $3e-4$. We train the model with a batch size of 16 for 40 epochs. Note that each video includes at least 1s audio or visual event, and 7202 video clips are annotated with more than one event category. We use 10,000 video clips with only video-level event labels for training. Following the official splits (Tian et al., 2020) of validation and test sets, we develop and test the model on the remaining 1879 videos with the segment-level annotations, *i.e.*, the speech event for audio starts at 1s and ends at 5s.

Audio-Visual Scene-Aware Dialog. In the audio-visual scene-aware dialog task, our model employs an advanced dialog generation framework that integrates audio and visual information to

1296 produce contextually relevant conversations. The dialog system utilizes a Transformer-based ar-
 1297 chitecture, which processes inputs from both modalities through separate encoders before merging
 1298 them in a fusion layer. This approach allows the model to understand the context provided by both
 1299 the audio and visual data streams effectively. The model is optimized using the Adam optimizer
 1300 with a learning rate of $1e - 4$ and a batch size of 64. Training is conducted for up to 30 epochs, with
 1301 early stopping based on performance on a validation set to prevent overfitting.

1302 **Audio-Visual Question-Answering.** For the AVQA task, our implementation focuses on integrating
 1303 spatial and temporal grounding techniques to accurately answer questions based on the video and
 1304 audio content. The system employs a dual-stream encoder that separately processes visual and audio
 1305 inputs. The encoded features are then combined using a co-attention mechanism that aligns audio
 1306 and visual elements relevant to the question context. This integration allows the model to focus on
 1307 specific segments of audio and video that are crucial for answering the given question. The model
 1308 is trained using the Adam (Kingma & Ba, 2014) optimizer with an initial learning rate of $3e - 4$,
 1309 reduced by a factor of 0.1 upon plateauing of validation loss. The system is trained for 40 epochs
 1310 with a batch size of 32.

1311 **Audio-Visual Segmentation.** For segmentation, we follow the prior work (Zhou et al., 2022), and
 1312 apply an upsampling decoder on features from the last self-attention layer to generate the final
 1313 segmentation mask. We use the binary cross entropy (BCE) loss between the prediction and ground-
 1314 truth masks for training. The models are trained for 20 epochs using the Adam optimizer (Kingma
 1315 & Ba, 2014) with a learning rate of $1e - 4$ and a batch size of 128.

1316 **Audio-Visual Source Separation.** For sound source separation, we follow the previous
 1317 method (Zhao et al., 2018; Mo & Morgado, 2023) and attach an audio U-Net decoder to our pre-
 1318 trained audio-visual encoders for separating sounds from the mixture. The decoder depth for self-
 1319 attention layers is 8, and the decoder receives the representations of the audio mixture and the visual
 1320 embeddings. We also apply multiple transposed convolutions and an output head to predict a time-
 1321 frequency separation mask. This separation mask is then used to multiply the input mixture STFT
 1322 to separate the audio. Similarly to (Zhao et al., 2018), the target masks refer to the time-frequency
 1323 bins where the source is the most dominant component in the mixture. The sound source separation
 1324 is achieved by optimizing a binary cross-entropy loss over these binary targets. The model is trained
 1325 for 20 epochs using the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $1e - 4$ and a
 1326 batch size of 128.

1327 **Dataset Details.** We evaluated our method using several prominent audio-visual datasets:

- 1328 • **Flick-SoundNet (Senocak et al., 2018):** a dataset consisting of natural soundscapes with
 1329 associated Flickr images with 4,500 audio-visual pairs for training and testing the model
 1330 on 250 audio-visual pairs of sounding objects and extended 250 non-sounding objects;
- 1331 • **VGG-Instruments (Hu et al., 2022):** contains video clips of musical instrument perfor-
 1332 mances, with 32k video clips of 10s lengths from 36 musical instrument classes, a subset
 1333 of VGG-Sound (Chen et al., 2020b), and each video only has one single instrument class;
- 1334 • **MUSIC (Zhao et al., 2018):** consists of 448 untrimmed YouTube music videos of solos
 1335 and duets from 11 instrument categories;
- 1336 • **VGG-Music (Mo & Morgado, 2023):** a dataset that features a collection of music videos
 1337 with annotations related to the genre and instruments present;
- 1338 • **VGGSound (Chen et al., 2020b):** a comprehensive dataset that includes a wide variety
 1339 of sound categories and corresponding visual scenes, which contains categories, such as
 1340 animals, instruments, vehicles, people, etc;
- 1341 • **AudioSet (Gemmeke et al., 2017):** a collection of 2,084,320 human-labeled 10-second
 1342 sound clips drawn from YouTube videos with 632 audio event classes;
- 1343 • **AVSBench (Zhou et al., 2022):** a benchmark for testing audio-visual synchronization and
 1344 alignment in diverse settings, including 4,932 videos (in total 10,852 frames) from 23 cat-
 1345 egories, including instruments, humans, animals, etc.
- 1346 • **MSR-VTT (Xu et al., 2016):** A large-scale video description dataset that includes 10,000
 1347 video clips, each paired with 20 human-annotated captions, useful for tasks involving video
 1348 understanding and retrieval.

- **LLP (Tian et al., 2020)**: The Look, Listen, and Parse (LLP) dataset contains densely labeled video segments that are used to train and evaluate models on tasks requiring fine-grained temporal understanding of video content.
- **MUSIC-AVQA (Li et al., 2022)**: A dataset specifically curated for audio-visual question answering in 11,849 YouTube video clips of 10 seconds long from 25 different event categories. It combines visual and audio clues to answer complex queries about the content and context of musical pieces.

Model Architecture and Training Details. We used a ViT-Base model as the backbone for all experiments, maintaining consistency with prior work to ensure comparability. The model incorporates a 12-layer transformer with 768 hidden dimensions and 12 attention heads. For processing audio-visual pairs, the audio input is represented as spectrograms with dimensions 128×128 , while the video input consists of frames resized to 224×224 . The model jointly processes these inputs through separate audio and visual encoders, followed by a cross-modal attention mechanism. Training was conducted for 100 epochs using the Adam optimizer with a learning rate of $1e - 4$ and a batch size of 128. The same configuration was applied across all downstream tasks, ensuring a uniform and fair experimental setup. Table 16 summarizes the key parameters of the model architecture and training pipeline.

Table 16: Model Configuration for Downstream Tasks.

Parameter	Value
Model Architecture	ViT-Base
Input Resolution	224x224
Batch Size	128
Optimizer	Adam
Learning Rate	1×10^{-4}
Training Epochs	50

G DATASET CONSTRUCTION AND QUALITY ANALYSIS

G.1 DATA CURATION WORKFLOW

To ensure precise and informative annotations, we developed a three-step multimodal captioning workflow that combines state-of-the-art tools for audio and video description with advanced aggregation techniques.

- **Video Captions:** For visual inputs, such as video frames, we use VideoLLaVA (Lin et al., 2023), a vision-language model specifically designed for video understanding. VideoLLaVA processes sampled frames from the video and generates multiple sentence-level captions that describe the visual content, including objects, actions, and scene attributes. The model’s ability to capture fine-grained visual details provides a strong foundation for multimodal representation.
- **Audio Captions:** For audio inputs, we utilize WavCaps (Mei et al., 2024), a model optimized for generating descriptive captions from audio signals. WavCaps processes the 10-second audio segments from each video, capturing key audio characteristics such as environmental sounds, speech, music, or other auditory cues. This step ensures that the audio component is accurately represented, complementing the visual descriptions.
- **Final Aggregation:** To produce a coherent and unified text description, we employ GPT-4 (OpenAI, 2023), a multimodal large language model. GPT-4 takes the outputs from VideoLLaVA and WavCaps as input and aggregates them into a single, semantically consistent caption. This final description integrates visual and auditory details into a cohesive narrative, ensuring cross-modal alignment and reducing redundancy. Additionally, GPT-4 is prompted to enhance the captions by resolving ambiguities and providing contextual information where necessary.

This structured workflow ensures that each sample in the dataset is annotated with high-quality, multimodal descriptions that capture the complementary nature of audio and video. By combining

specialized models for each modality with an advanced language model for integration, our approach achieves detailed and contextually rich annotations suitable for a wide range of downstream tasks.

G.2 ALIGNMENT AND FILTERING

Our alignment and filtering pipeline is a critical step in ensuring the high quality of the *ACAV-1M* dataset, particularly in improving cross-modal alignment between audio, video, and text. We utilize ImageBind, a powerful representation learning framework, to evaluate and optimize the alignment quality during the curation process. The alignment quality is measured using cosine similarity scores across three modalities: audio-visual, text-audio, and text-visual. Table 17 provides a comparative analysis of the alignment scores for the unfiltered *ACAV-100M* dataset and the curated *ACAV-1M* dataset. The results demonstrate substantial improvements in alignment quality across all modalities after filtering. The improvements in alignment scores demonstrate that our filtering pipeline significantly enhances the dataset’s quality. By selectively retaining samples with high alignment scores, we ensure that *ACAV-1M* provides clean, synchronized audio-visual pairs and detailed, semantically consistent captions. These properties are essential for training robust multimodal models and achieving superior performance across a wide range of downstream tasks. The filtering process effectively addresses the inherent noise and misalignments in the unfiltered *ACAV-100M* dataset, ensuring that *ACAV-1M* is well-suited for applications requiring fine-grained audio-visual-textual integration.

Table 17: Alignment quality comparison.

Dataset	Audio-Visual Alignment	Text-Audio Alignment	Text-Visual Alignment
<i>ACAV-100M</i>	0.42	0.38	0.44
<i>ACAV-1M</i>	0.62	0.58	0.65

G.3 DATASET COMPOSITION AND DISTRIBUTION

Category Distribution. To evaluate the diversity and representativeness of the *ACAV-1M* dataset, we conducted an analysis of the category distribution across its samples. This analysis ensures that the dataset provides a balanced representation of common audio-visual scenarios, enabling broad applicability across various downstream tasks. Table 18 presents the proportion of six major categories within the dataset. Music, accounting for 30% of the dataset, represents the largest category, reflecting the prevalence of music-related content in audio-visual datasets. Nature sounds and scenes constitute 20%, emphasizing environmental diversity. Speech/Dialogues make up 15%, capturing conversational and narrative contexts critical for tasks like audio-visual question answering and scene-aware dialogue generation. Vehicles and sports each account for 10%, covering dynamic content that often involves synchronized motion and sound. Finally, the “Others” category, comprising 15%, includes diverse scenarios such as industrial environments, urban landscapes, and artistic performances. The balanced distribution of categories ensures that the dataset supports the development of generalizable models while catering to specialized applications. By avoiding overrepresentation of any single category, *ACAV-1M* provides a robust foundation for training multimodal systems capable of handling diverse real-world scenarios. This careful curation strengthens the dataset’s utility for benchmarking and improving state-of-the-art audio-visual learning methods.

Table 18: Dataset composition by category in *ACAV-1M*.

Category	Proportion (%)
Music	30
Nature	20
Speech/Dialogues	15
Vehicles	10
Sports	10
Others	15

Table 19: Temporal distribution of videos in *ACAV-1M*.

Year	Number of Videos (%)
2010-2012	10
2013-2015	15
2016-2018	20
2019-2020	25
2021	30

Temporal Distribution. To ensure the temporal representativeness of the *ACAV-1M* dataset, we analyzed the year-wise distribution of video samples. This analysis helps verify that the dataset spans a wide range of time periods, reflecting changes in audio-visual content trends and maintaining relevance across various contexts. Table 19 summarizes the temporal distribution of videos in *ACAV-1M*. The dataset includes content from 2010 to 2023, with a noticeable peak in data collection during 2021, accounting for 30% of the total dataset. The years 2019-2020 contribute 25% of the data, followed by 20% from 2016-2018, 15% from 2013-2015, and 10% from the earliest period, 2010-2012. This temporal distribution ensures that *ACAV-1M* captures diverse audio-visual patterns and contexts, enabling models trained on it to generalize effectively across different timeframes. The peak in 2021 likely reflects increased availability of high-quality audio-visual content during this period. By including data from over a decade, the dataset provides a rich temporal context, enhancing its utility for time-sensitive applications and historical trend analysis in audio-visual learning.

Data Quality Control. To ensure the high quality and reliability of the *ACAV-1M* dataset, we implemented a rigorous quality control process. This process involved a random inspection of 10,000 samples to evaluate their adherence to the alignment criteria established during the curation workflow. The alignment criteria include semantic consistency across audio, video, and text modalities, as well as temporal synchronization between audio and visual streams. The inspection results showed that approximately 98.4% of the evaluated samples met the alignment criteria, demonstrating the effectiveness of our automated filtering pipeline. The remaining 1.6% of samples, which exhibited issues such as temporal misalignment or semantic inconsistencies, were flagged for manual review. These flagged samples were either corrected or removed to maintain the dataset’s integrity. This stringent quality control process ensures that *ACAV-1M* provides clean, well-aligned data, minimizing noise and errors that could impact downstream tasks. By combining automated filtering with manual inspection, we achieve a robust dataset that serves as a reliable audio-visual benchmark.

G.4 EFFECTIVENESS OF FILTERING PIPELINE

Filtered vs. Unfiltered Dataset. To evaluate the impact of our filtering pipeline, we compared the performance of models trained on the filtered *ACAV-1M* dataset against those trained on the unfiltered *ACAV-100M* dataset. This comparison highlights the advantages of using *ACAV-1M* in terms of improved alignment quality and enhanced downstream performance. Table 20 summarizes the results across three key tasks: audio-visual classification, audio-visual retrieval, and source localization. For audio-visual classification, models trained on *ACAV-1M* achieved an accuracy of 68.23%, representing a 6.78 percentage point improvement over the unfiltered dataset. Similarly, in audio-visual retrieval, the Recall@1 (R@1) score increased from 22.56% to 26.57%, demonstrating a significant enhancement in retrieval precision. For source localization, precision improved by 8.44 percentage points, from 50.23% to 58.67%. These results indicate that the filtering process in *ACAV-1M* successfully enhances the dataset’s quality, leading to substantial performance gains across tasks. The improvements highlight the importance of alignment and filtering in reducing noise and improving data coherence, ultimately enabling more robust and accurate model training. This comparison underscores the necessity of curating high-quality datasets for *ACAV-1M*.

Data Loss During Filtering. To achieve the high alignment quality of *ACAV-1M*, a multi-stage filtering process was implemented. This process ensures that only samples meeting strict alignment criteria across language, instance, and temporal dimensions are included in the final dataset. Table 21 details the proportion of data excluded at each filtering stage. The largest reduction occurred during the temporal alignment filtering stage, which excluded 38% of the total samples, followed by instance alignment filtering (36%) and language alignment filtering (25%). These stages collectively

Table 20: Filtered vs. Unfiltered dataset performance.

Task	Metric	ACAV-100M (Unfiltered)	ACAV-1M (Filtered)
Audio-Visual Classification	Accuracy (%)	61.45	68.23 (+6.78)
Audio-Visual Retrieval	R@1 (%)	22.56	26.57 (+4.01)
Source Localization	Precision (%)	50.23	58.67 (+8.44)

resulted in the exclusion of 99% of the original dataset samples, highlighting the stringent quality requirements applied during curation. While this level of exclusion significantly reduces the dataset’s size, it ensures that the remaining samples exhibit high-quality alignment across modalities. The filtered dataset prioritizes precision and coherence, which are critical for robust training and evaluation in audio-visual learning. This rigorous filtering process is essential for eliminating noise and inconsistencies, ultimately enabling the development of models that achieve superior performance on downstream tasks.

Table 21: Data exclusion across filtering stages.

Filtering Stage	Proportion Excluded (%)
Language Alignment Filtering	25
Instance Alignment Filtering	36
Temporal Alignment Filtering	38
Total Excluded Data	99

Optimized Thresholds. To improve the effectiveness of the filtering process, we independently optimized the similarity thresholds for different alignment types: instance alignment, temporal alignment, and language alignment. This approach avoids the limitations of a uniform threshold and tailors the alignment process to the specific characteristics of each modality. Table 22 reports the optimized thresholds and the corresponding performance improvements compared to using a uniform 50% threshold across all alignment types. For instance alignment, the optimized threshold remained at 50%, resulting in a 3.24% performance improvement. Temporal alignment achieved the best results with a threshold of 70%, yielding a 4.12% improvement. Language alignment benefited from a threshold of 60%, with the largest performance gain of 5.13%. These results demonstrate that fine-tuning thresholds for each alignment type enhance the overall quality of the filtered dataset. The optimized thresholds ensure that samples meeting the criteria are retained, improving data alignment while maintaining diversity. This nuanced approach strengthens the dataset’s ability to support diverse downstream tasks, ultimately leading to better model performance.

Table 22: Impact of optimized similarity thresholds.

Alignment Type	Threshold (%)	Performance Improvement (%)
Instance Alignment	50	+3.24
Temporal Alignment	70	+4.12
Language Alignment	60	+5.13

G.5 EVALUATION OF IMAGEBIND AND CAPTIONING QUALITY

ImageBind Cross-Modal Alignment. We conducted an evaluation of ImageBind’s cross-modal alignment accuracy on a representative subset of 10,000 samples from the dataset. The goal of this evaluation was to assess ImageBind’s ability to align audio, video, and text modalities effectively, a critical aspect of our dataset curation process. Table 23 summarizes the results of the evaluation. ImageBind achieved an alignment accuracy of 88.45%, demonstrating its reliability in ensuring high-quality alignment across modalities. This accuracy indicates that ImageBind provides a robust foundation for filtering and aligning samples in *ACAV-1M*, significantly contributing to the dataset’s overall quality. The strong performance of ImageBind validates its use as a key component in our alignment and filtering pipeline. By leveraging its cross-modal capabilities, we ensure that *ACAV-1M*

offers highly aligned and semantically coherent samples suitable for a wide range of downstream tasks. This evaluation underscores the importance of reliable cross-modal alignment tools in creating high-quality multimodal datasets.

Table 23: Cross-modal alignment accuracy of ImageBind.

Metric	Accuracy (%)
Cross-Modal Alignment Accuracy	88.45

Human Evaluation of Captions. To assess the quality of GPT-4-generated captions in *ACAV-1M*, we conducted a human evaluation on a subset of the dataset. Participants rated the captions for two key metrics: relevance, which measures how well the captions describe the corresponding audio-visual content, and correctness, which assesses the factual accuracy of the captions. Table 24 presents the results of this evaluation. GPT-4-generated captions scored 92.34% in relevance, demonstrating their strong alignment with the underlying content. The correctness score of 85.12% indicates that the captions are generally accurate but may require minor improvements for specific edge cases. These results highlight the reliability of GPT-4 in generating high-quality captions that are both contextually relevant and semantically accurate. By using GPT-4 to aggregate and refine captions from audio and video inputs, we ensure that *ACAV-1M* provides detailed and coherent text descriptions. This enhances the dataset’s utility for tasks such as audio-visual retrieval, question answering, and scene-aware dialogue.

Table 24: Human evaluation of captions.

Metric	Relevance (%)	Correctness (%)
GPT-4-Generated Captions	92.34	85.12

Reliability of ImageBind Filtering. While ImageBind was not explicitly trained on text-audio pairs, its generalization capabilities make it a critical component of our filtering pipeline. To ensure its reliability for cross-modal alignment tasks, we conducted a detailed evaluation of its performance across audio-visual and text-based pairs. Table 25 summarizes the results of this validation. ImageBind achieved high alignment accuracy for audio-visual (88.45%) and text-visual (89.67%) pairs, demonstrating strong reliability in these modalities. For text-audio pairs, ImageBind attained a slightly lower accuracy of 85.12%, reflecting moderate reliability. While its performance on text-audio pairs is not as robust as for other alignments, it is still sufficient for maintaining high-quality filtering in *ACAV-1M*. These results affirm the utility of ImageBind in our data curation process. Its capability to generalize across modalities ensures that the dataset maintains a high degree of alignment and semantic coherence, even in cases where direct training data is limited. This evaluation supports the robustness of ImageBind as a filtering tool for building multimodal datasets.

Table 25: Reliability of ImageBind for cross-modal filtering.

Alignment Type	Accuracy (%)	Comment
Audio-Visual	88.45	Highly Reliable
Text-Audio	85.12	Moderate Reliability
Text-Visual	89.67	Highly Reliable

G.6 DATASET QUALITY CONTROL

Quality Control Process. To ensure alignment consistency in *ACAV-1M*, a comprehensive quality control process was implemented. This process involved a detailed manual inspection of a randomly selected subset of samples and the establishment of procedures for handling misalignments identified after dataset curation. Table 26 summarizes the results of the manual inspection. A total of 10% of the dataset was manually inspected, with an alignment accuracy of 98.4%. Only 1.6% of the inspected samples were found to be misaligned. These misaligned samples were flagged for manual

review and either corrected or excluded from the dataset to maintain high-quality standards. In addition to the initial inspection, a post-curation process was established to address misalignments identified during downstream task evaluations. Any flagged samples undergo a thorough review and are updated as necessary to ensure the dataset remains reliable for diverse applications. This ongoing quality assurance process reinforces the integrity of *ACAV-1M*, ensuring that it serves as a dependable resource for audio-visual learning tasks.

Table 26: Manual quality inspection of *ACAV-1M*.

Metric	Percentage (%)
Manually Inspected Samples	10.0
Alignment Accuracy	98.4
Misaligned Samples	1.6

G.7 EFFECT OF ALIGNMENT ON PERFORMANCE

ACAV-100M vs. *ACAV-1M*. To assess the impact of our filtering and alignment process, we compared the performance of models trained on the unfiltered *ACAV-100M* dataset with those trained on the filtered *ACAV-1M* dataset. This evaluation highlights the effectiveness of the curated dataset in improving task performance by enhancing data quality and alignment. Table 27 shows the results of this comparison across three key downstream tasks. For audio-visual classification, models trained on *ACAV-1M* achieved an accuracy of 68.23%, outperforming the unfiltered dataset by 4.78 percentage points. Similarly, in audio-visual retrieval, the Recall@1 (R@1) score improved by 4.23 percentage points, demonstrating the benefits of enhanced alignment and semantic coherence. In source localization, the precision increased by 5.82 percentage points, further emphasizing the impact of high-quality alignment in *ACAV-1M*. These results validate the importance of the filtering and alignment processes implemented in *ACAV-1M*. By prioritizing data quality and coherence, *ACAV-1M* enables models to achieve superior performance across a variety of audio-visual tasks, confirming its value as a robust benchmark dataset.

Table 27: Performance comparison between *ACAV-100M* and *ACAV-1M*.

Task	Metric	ACAV-100M (Unfiltered)	ACAV-1M (Filtered)
Audio-Visual Classification	Accuracy (%)	63.45	68.23 (+4.78)
Audio-Visual Retrieval	R@1 (%)	22.34	26.57 (+4.23)
Source Localization	Precision (%)	52.85	58.67 (+5.82)

Similarity Distribution. The filtering process in *ACAV-1M* was designed to improve the alignment quality across audio-visual, text-audio, and text-visual pairs. To evaluate its impact, we analyzed the similarity distributions before and after filtering. Table 28 presents the mean similarity scores for each pair type in the unfiltered *ACAV-100M* dataset and the filtered *ACAV-1M* dataset, alongside the thresholds applied during filtering. For audio-visual pairs, the mean similarity increased from 0.42 in *ACAV-100M* to 0.62 in *ACAV-1M*, demonstrating a significant improvement in alignment quality. Similar trends were observed for text-audio pairs (from 0.38 to 0.58) and text-visual pairs (from 0.44 to 0.65). The thresholds applied during filtering (set at 0.50 for all similarity types) ensured that only samples meeting a high alignment standard were retained in *ACAV-1M*. These results highlight the effectiveness of the filtering process in enhancing the alignment quality across modalities. By enforcing rigorous thresholds, *ACAV-1M* provides a more coherent and semantically consistent dataset, which is essential for training robust audio-visual models and achieving superior performance across diverse downstream tasks.

Statistical Similarity Comparison. To provide a detailed evaluation of alignment quality, we conducted a statistical analysis of similarity scores before and after filtering. Table 29 presents the mean similarity scores for audio-visual, text-audio, and text-visual pairs in the pre-filtered and post-filtered datasets, along with the percentage improvement achieved through the filtering process. For audio-visual pairs, the mean similarity increased from 0.42 (pre-filtering) to 0.62 (post-filtering), representing a 47.62% improvement. Text-audio pairs showed an even greater increase, with mean

Table 28: Similarity distribution before and after filtering.

Similarity Type	ACAV-100M (Mean)	ACAV-1M (Mean)	Threshold Applied
Audio-Visual	0.42	0.62	0.50
Text-Audio	0.38	0.58	0.50
Text-Visual	0.44	0.65	0.50

similarity scores improving by 52.63% (from 0.38 to 0.58). Similarly, text-visual pairs exhibited a 47.73% improvement, with mean scores increasing from 0.44 to 0.65. These results demonstrate the effectiveness of the filtering pipeline in enhancing the alignment quality across modalities. By selecting high-quality samples based on rigorous similarity thresholds, the *ACAV-1M* dataset achieves significantly better alignment compared to the unfiltered *ACAV-100M* dataset. This improvement ensures a more reliable and coherent dataset for training multimodal models and achieving superior performance across downstream tasks.

Table 29: Statistical similarity comparison before and after filtering.

Similarity Type	Mean (Pre-Filtering)	Mean (Post-Filtering)	Improvement (%)
Audio-Visual	0.42	0.62	+47.62
Text-Audio	0.38	0.58	+52.63
Text-Visual	0.44	0.65	+47.73

G.8 VIDEO AND AUDIO SEGMENTATION DETAILS

Segment Cutting and Semantic Integrity. To ensure that segment cutting does not compromise the quality of the *ACAV-1M* dataset, we evaluated the preservation of semantic alignment and temporal consistency after segmenting videos into 10-second clips. Table 30 presents the alignment scores, semantic relevance, and temporal synchronization metrics both before and after segmentation, along with the percentage change. The results demonstrate minimal degradation in alignment quality and semantic integrity due to segment cutting. The audio-visual alignment score decreased slightly by 3.17%, from 0.63 to 0.61, indicating that the segmentation process maintains most of the alignment between modalities. Semantic relevance showed a small decline of 2.44%, reflecting the retention of meaningful context in the segmented clips. Temporal synchronization exhibited the least change, with a marginal reduction of 1.39%, further supporting the temporal consistency of the segments. These findings confirm that segment cutting into 10-second clips introduces negligible impact on the dataset’s semantic and alignment quality. This ensures that the *ACAV-1M* dataset provides high-quality training samples suitable for a wide range of audio-visual learning tasks.

Table 30: Alignment and semantic integrity after segment cutting.

Metric	Pre-Segmentation	Post-Segmentation	Change (%)
Audio-Visual Alignment Score	0.63	0.61	-3.17
Semantic Relevance Score	0.82	0.80	-2.44
Temporal Synchronization	0.72	0.71	-1.39

Semantic Integrity of 10-second Segments. To ensure that the segmentation process preserves semantic coherence, we employed logical boundaries, such as scene changes or natural pauses, as cut points when dividing videos into 10-second segments. This approach minimizes disruptions to the narrative or contextual flow of the video and audio content. Semantic integrity was further validated by evaluating 10,000 randomly selected segments from the dataset. As shown in Table 31, 95.2% of the segments retained semantic coherence, indicating that the segmentation process successfully preserved the contextual and thematic consistency of the content. Only 4.8% of the segments were identified as misaligned or semantically inconsistent. These results confirm that the segmentation process in *ACAV-1M* maintains a high standard of semantic integrity, ensuring the dataset’s suit-

ability for downstream tasks that rely on coherent and contextually relevant data. This evaluation underscores the robustness of the curation pipeline in producing high-quality training samples.

Table 31: Semantic integrity evaluation.

Metric	Percentage (%)
Semantically Coherent Segments	95.2
Misaligned Segments	4.8

H MORE EXPERIMENTAL ANALYSIS

H.1 COMPARATIVE ANALYSIS WITH AUDIOSET

Performance Comparison with AudioSet. To evaluate the advantages of *ACAV-1M* over AudioSet (Gemmeke et al., 2017), we conducted a performance comparison across multiple downstream tasks. The results, summarized in Table 32, clearly demonstrate the benefits of *ACAV-1M*’s alignment-focused curation pipeline. For audio-visual classification, models trained on *ACAV-1M* achieved an accuracy of 68.23%, outperforming AudioSet by 2.51%. This improvement highlights the value of *ACAV-1M*’s precise alignment in enhancing model performance on classification tasks. In the audio-visual retrieval task, *ACAV-1M* achieved a Recall@1 (R@1) score of 26.57%, surpassing AudioSet by 1.72%. This indicates that *ACAV-1M*’s curation pipeline leads to more semantically aligned and contextually rich data, which is critical for retrieval tasks. For source localization, *ACAV-1M* achieved a precision of 58.67%, significantly exceeding AudioSet’s 54.32% by 4.35%. This highlights the effectiveness of *ACAV-1M*’s temporal and instance alignment filtering in producing data that supports fine-grained localization tasks. These results underline *ACAV-1M*’s superior alignment and curation quality, demonstrating its potential to drive advancements in audio-visual representation learning across various tasks.

Table 32: Comparison between *ACAV-1M* and AudioSet on downstream tasks.

Task	Metric	AudioSet	<i>ACAV-1M</i>
Audio-Visual Classification	Accuracy (%)	65.72	68.23 (+2.51)
Audio-Visual Retrieval	R@1 (%)	24.85	26.57 (+1.72)
Source Localization	Precision (%)	54.32	58.67 (+4.35)

Task Support Comparison. We conducted a side-by-side analysis of the tasks supported by *ACAV-1M* and AudioSet to highlight the unique capabilities of our dataset. Table 33 presents this comparison. While both datasets support foundational tasks such as audio-visual classification, retrieval, and source localization, *ACAV-1M* extends its utility to advanced tasks like temporal segmentation, scene-aware dialogue, and audio-visual question-answering. The inclusion of these advanced tasks underscores the versatility of *ACAV-1M*, made possible through its robust curation pipeline that emphasizes alignment and semantic richness. Temporal segmentation benefits from *ACAV-1M*’s focus on temporal consistency, while scene-aware dialogue and question-answering tasks leverage the dataset’s multimodal captioning and fine-grained alignment. This expanded task support makes *ACAV-1M* a more comprehensive resource for multimodal learning and related applications.

Table 33: Task support comparison between AudioSet and *ACAV-1M*.

Task	AudioSet	<i>ACAV-1M</i>
Audio-Visual Classification	✓	✓
Audio-Visual Retrieval	✓	✓
Source Localization	✓	✓
Temporal Segmentation	✗	✓
Scene-Aware Dialogue	✗	✓
Audio-Visual Question-Answering	✗	✓

1782 H.2 ADVANCED BASELINES: IMAGEBIND IN RETRIEVAL

1783
1784 **ImageBind Baseline for Audio-Visual Retrieval.** We evaluated ImageBind as a baseline for
1785 audio-visual retrieval and compared it with a model trained on the *ACAV-1M* dataset. Table 34
1786 illustrates the performance metrics across recall@1, recall@5, and recall@10. The results demon-
1787 strate that the *ACAV-1M*-trained model consistently outperforms ImageBind, with notable improve-
1788 ments in all retrieval metrics. These findings validate the effectiveness of *ACAV-1M* in training robust
1789 audio-visual retrieval models, highlighting the benefits of its high-quality alignment and compre-
1790 hensive curation pipeline. The gains across higher recall thresholds (R@5 and R@10) emphasize the
1791 dataset’s capability to improve retrieval precision and robustness in challenging multimodal tasks.

1792 Table 34: Audio-Visual retrieval performance: ImageBind vs. model trained on *ACAV-1M*.

Method	R@1 (%)	R@5 (%)	R@10 (%)
ImageBind	22.34	51.25	62.84
<i>ACAV-1M</i> Model	26.57	58.78	70.26

1798 H.3 ALTERNATIVE MODELS FOR ALIGNMENT: FREEBIND AND OMNIBIND

1800
1801 **Comparative Evaluation with FreeBind and OmniBind.** We evaluated the performance of
1802 *ACAV-1M* using advanced alignment models FreeBind (Wang et al., 2024a) and OmniBind (Wang
1803 et al., 2024b), comparing them with ImageBind as the baseline. Table 35 presents the results, high-
1804 lighting the cross-modal alignment accuracy, temporal alignment accuracy, and overall task perfor-
1805 mance improvements achieved with each model. The results indicate that OmniBind outperforms
1806 both ImageBind and FreeBind across all metrics, achieving a 92.18% cross-modal alignment accu-
1807 racy and 89.63% temporal alignment accuracy, leading to a +7.30 improvement in downstream task
1808 performance. FreeBind also shows significant gains over the baseline, with a +2.52 performance
1809 increase. These findings demonstrate the potential of leveraging advanced alignment models to fur-
1810 ther enhance the utility of *ACAV-1M* in multimodal learning tasks, particularly in scenarios requiring
1811 precise temporal and semantic alignment.

1812 Table 35: Comparison of alignment models on *ACAV-1M*.

Method	Cross-Modal Alignment (%)	Temporal Alignment (%)	Task Results (%)
ImageBind	88.45	85.12	68.23
FreeBind	91.34	88.25	70.75
OmniBind	92.18	89.63	75.53

1818 H.4 EVALUATION WITH AUDIOLLM FOR CAPTIONING

1820
1821 **Comparison Between VideoLLaVA and AudioLLM.** To assess the effectiveness of using dif-
1822 ferent captioning models in our pipeline, we compared VideoLLaVA, a multimodal model, with
1823 AudioLLM (WavCaps (Mei et al., 2024)) for audio captioning. Table 36 provides a detailed com-
1824 parison of model performance across downstream tasks when trained with captions generated by
1825 these systems. The results show that VideoLLaVA achieves superior performance across all tasks,
1826 with an accuracy of 68.23% in audio-visual classification, 26.57% in R@1 for audio-visual retrieval,
1827 and a precision of 58.67% in source localization. In comparison, AudioLLM performs competitively
1828 but falls short, demonstrating the advantages of using a model with multimodal capabilities. These
1829 findings highlight the importance of integrating multimodal alignment in the captioning process to
1830 enhance task-specific performance, affirming our choice of VideoLLaVA for generating captions in
1831 the *ACAV-1M* pipeline.

1832 H.5 RETRAINING IMAGEBIND WITH *ACAV-1M*

1833
1834 **Impact of Retraining ImageBind on *ACAV-1M*.** To evaluate the impact of *ACAV-1M* on foun-
1835 dational models, we retrained ImageBind using our dataset. Table 37 compares the performance
of the original ImageBind model and the retrained version on several cross-modal tasks, including

Table 36: Comparison between VideoLLaVA and AudioLLM for captioning.

Task	Metric	VideoLLaVA	AudioLLM
Audio-Visual Classification	Accuracy (%)	68.23	66.85
Audio-Visual Retrieval	R@1 (%)	26.57	25.34
Source Localization	Precision (%)	58.67	57.12

audio-visual classification, retrieval, and source localization. The results demonstrate substantial improvements across all tasks after retraining ImageBind on *ACAV-1M*. For audio-visual classification, the retrained model achieved an accuracy of 69.12%, a 3.4% improvement over the original. In audio-visual retrieval, the R@1 metric increased significantly from 22.34% to 27.45%, and source localization precision rose from 54.32% to 59.88%. These findings underscore the utility of *ACAV-1M* in enhancing foundational models, particularly in tasks requiring fine-grained multimodal understanding. Retraining ImageBind on *ACAV-1M* not only improves its baseline performance but also validates the dataset’s alignment-focused curation pipeline as a valuable resource for advancing cross-modal learning.

Table 37: Performance of original and retrained ImageBind models.

Task	Metric	Original ImageBind	Retrained ImageBind
Audio-Visual Classification	Accuracy (%)	65.72	69.12
Audio-Visual Retrieval	R@1 (%)	22.34	27.45
Source Localization	Precision (%)	54.32	59.88

H.6 CONTENT-BASED FILTERING PARAMETERS

Adjusting Filtering for Different Content Types. To further improve task-specific performance, we optimized filtering parameters for different content types, including music, ambient sounds, and narration. These adjustments were based on the unique alignment and synchronization requirements of each modality. Table 38 summarizes the performance improvements achieved through this content-specific filtering approach. The results indicate that tailored filtering significantly enhances performance across tasks. For instance, optimized filtering for music content led to a 2.62% improvement in source separation accuracy, while adjustments for ambient sounds improved audio-visual parsing by 1.73%. Similarly, retrieval tasks for narration saw a 1.55% increase in performance. These findings demonstrate the value of adapting alignment requirements to the characteristics of specific content types, highlighting the flexibility and scalability of the *ACAV-1M* curation pipeline in addressing diverse multimodal tasks.

Table 38: Impact of Content-Specific Filtering Parameters on Task Performance.

Content Type	Task	Baseline (%)	Optimized Filtering (%)
Music	Source Separation	65.72	68.34
Ambient Sounds	Audio-Visual Parsing	70.45	72.18
Narration	Retrieval	26.57	28.12

H.7 BALANCING STRICTNESS AND DIVERSITY IN FILTERING

Analysis of Data Loss and Diversity. To ensure a balanced approach between strictness and data diversity during filtering, we evaluated the impact of varying similarity thresholds on the dataset composition and downstream task performance. Table 39 highlights the trade-offs observed for thresholds of 25%, 50%, and 75%. A threshold of 50% emerged as the optimal setting, retaining 90% of the original dataset while preserving 96% of task diversity. This threshold also achieved the highest classification accuracy of 68.23%, reflecting the effectiveness of this setting in maintaining alignment quality without excessively reducing data diversity. In contrast, a stricter threshold of 75% reduced the dataset size to 80%, leading to a slight drop in task diversity and a notable decrease

in classification accuracy to 66.23%. On the other hand, a more lenient threshold of 25% retained 95% of the data but resulted in reduced alignment quality, as reflected in the lower classification accuracy of 65.45%. These findings underscore the importance of selecting appropriate thresholds to balance data quality, quantity, and task-specific diversity, demonstrating the flexibility of the *ACAV-1M* curation pipeline to optimize for different application scenarios.

Table 39: Impact of Filtering Thresholds on Data Loss and Diversity.

Threshold (%)	Data Retained (%)	Task Diversity (%)	Classification Accuracy (%)
75	80	92	66.23
50	90	96	68.23
25	95	97	65.45

H.8 BIAS DETECTION AND MITIGATION IN CAPTIONING

Addressing Biases in Large Language Models. To ensure fairness and representativeness in captions generated by large language models (LLMs), we conducted a detailed evaluation of demographic, regional, and content-specific biases. Using a subset of 10,000 captions, we identified areas of bias and implemented mitigation techniques, including prompt engineering and manual review. These efforts focused on reducing unintended biases while maintaining the semantic integrity of the captions. Table 40 summarizes the bias scores before and after mitigation. The results demonstrate significant reductions in all evaluated categories. Demographic bias, for example, decreased by 43.75%, while regional and content-specific biases were reduced by 41.38% and 36.00%, respectively. This confirms the effectiveness of the implemented bias mitigation strategies. By addressing these biases, the captions in *ACAV-1M* provide a fairer representation across demographic, regional, and content dimensions, supporting ethical and inclusive dataset usage. This work also highlights the importance of proactive bias detection and mitigation in multimodal datasets.

Table 40: Bias Evaluation and Mitigation in Captions.

Bias Type	Baseline Bias Score	Post-Mitigation Bias Score	Reduction (%)
Demographic Bias	0.32	0.18	43.75
Regional Bias	0.29	0.17	41.38
Content-Specific Bias	0.25	0.16	36.00

H.9 PERFORMANCE BOTTLENECKS IN SPECIFIC SCENARIOS

Long Video Processing and Noisy Environments. The performance of *ACAV-1M* was evaluated in two challenging scenarios: long video processing and noisy environments, as reported in Table 41. For long video processing, the baseline classification accuracy was observed at 63.45%, indicating a need for methods capable of maintaining semantic coherence and temporal alignment over extended durations. A proposed solution is the use of temporal chunking, which segments long videos into manageable clips while preserving their semantic continuity. This approach allows models to process each chunk independently, reducing computational overhead and mitigating the risk of performance degradation over time. In noisy environments, the source separation task achieved a baseline accuracy of 58.12%, reflecting the challenges posed by environmental noise interfering with audio-visual alignment. To address this, noise-robust models, such as those incorporating advanced denoising techniques or using noise-augmented training data, are recommended. These models aim to enhance resilience to background interference, improving alignment accuracy and task performance.

Table 41: Performance Bottlenecks in Challenging Scenarios.

Scenario	Task	Baseline Accuracy (%)	Proposed Solutions
Long Video Processing	Classification	63.45	Temporal Chunking
Noisy Environments	Source Separation	58.12	Noise Robust Models